



Research Institute
Social Cohesion

RDC

SOEP

SCP Documentation

German Social Cohesion Panel

SCP 2021-22 W1-2 Codebook HHRF: Weights for Households (German)



German Social Cohesion Panel

Established in 2021, the German Social Cohesion Panel (SCP) is a wide-ranging representative longitudinal study of private households in Germany, carried out in collaboration of the Research Institute Social Cohesion (RISC) and the German Socio-Economic Panel (SOEP).

The aim of the SCP Documentation is to thoroughly document the survey's data collection and data processing.

Recommended Citation

Groh-Samberg, O., Axenfeld, J. B., Gerlitz, J.-Y., Cornesse, C., Kroh, M., Lengfeld, H., Liebig, S., Minkus, L., Reinecke, J., Teichler, N., Traunmüller, R., & Zinn, S. (2026). SCP 2021-22 W1-2 - Codebook HHRF: Weights for Households (German). *German Social Cohesion Panel 2021-2022 - Wave 1-2*. Bremen and Berlin: RDC-RISC/SOCIUM, SOEP/DIW Berlin. doi:10.60532/scp.2021-22.w1-2.v1

- ▶ **Authors:** Olaf Groh-Samberg, Julian B. Axenfeld, Jean-Yves Gerlitz, Carina Cornesse, Martin Kroh, Holger Lengfeld, Stefan Liebig, Lara Minkus, Jost Reinecke, Nils Teichler, Richard Traunmüller, Sabine Zinn

- ▶ **Contributors:** Cosima Adams, Anton Bochert, Martin Gerike, Josefine Kuhrmeier, Anna-Tabea Müller, Eric Nissen, Sebastian Rueda-Uribe, Rainer Siegers, Hans Walter Steinhauer, Knut Wenzig, Julia Witton (Project Members), infas (Data Collector)

- ▶ **Publisher:** RDC-RISC
SOCIUM, University of Bremen
P.O. Box 330 440
28334 Bremen
Germany

SOEP
DIW Berlin
German Socio-Economic Panel (SOEP)
Mohrenstr. 58
10117 Berlin
Germany

- ▶ **DOI:** 10.60532/scp.2021-22.w1-2.v1

- ▶ **Website:** www.fgz-risc-data.de
www.diw.de



The text of this publication is published under the Creative Commons license CC BY-SA 4.0 Attribution-ShareAlike 4.0 International. The exact wording of the license CC BY-SA 4.0 can be found here:

<https://creativecommons.org/licenses/by-sa/4.0/>

SCP Documentation

German Social Cohesion Panel

SCP 2021-22 W1-2 Codebook HHRF: Weights for Households (German)

Inhaltsverzeichnis

1	Allgemeine Informationen	2
2	Identifikatoren	2
	hid – Aktuelle Haushalts-ID	2
	cid – Ursprüngliche Haushalts-ID	3
3	Befragungskontext	3
	wave – Erhebungswelle	3
4	Statistische Gewichtungsfaktoren	4
	design – Inverse Ziehungswahrscheinlichkeit	4
	hhrf – Hochrechnungsfaktor	5
5	Inverse Bleibewahrscheinlichkeit	6
	hbleib – Inverse Bleibewahrscheinlichkeit	6

1 Allgemeine Informationen

Der HHRF-Datensatz enthält Gewichtungsfaktoren für die Haushalte im SCP. Jeder Haushalt (HID), der in einer bestimmten Erhebungswelle (WAVE) auf den Haushaltsfragebogen geantwortet hat, hat eine Zeile im Datensatz. Haushalte, deren Ankerperson nicht an der Erhebungswelle teilgenommen hat oder die Befragung vor Beginn des Haushaltsfragebogens abgebrochen hat, sind nicht im Datensatz enthalten.

An einigen Stellen in der Dokumentation und in den Daten werden Jahreszahlen z. B. für die Bezeichnung von Variablen und des Fragebogeninstruments verwendet. Diese Jahreszahlen orientieren sich stets am Feldstart der Datenerhebung der entsprechenden Erhebungswelle.

2 Identifikatoren

hid - Aktuelle Haushalts-ID

21100003	3
21100009	3
21100010	1
21100012	1
21100016	3
21100020	3
21100021	3
21100032	1
21100037	1
21100039	3
21100044	1
21100045	1
21100049	3
21100050	1
21100058	1
... (13029 Zeilen unterdrückt)	25800
21137972	1
21137973	1
21137976	1
21137978	2
21137979	3
21137985	2
21137987	1
21137991	3
21138000	1
22103378	1
22103896	1
22115150	1
22117540	1
22119085	1
22125622	1

Diese Kennziffer gruppiert Individuen in ihre zugehörigen Haushalte zum Zeitpunkt der aktuellen Erhebungswelle. Das bedeutet, dass die HID einer Person sich über die Zeit verän-

dern kann, zum Beispiel wenn ein erwachsenes Kind aus dem elterlichen Haushalt auszieht und einen eigenen Haushalt eröffnet.

cid – Ursprüngliche Haushalts-ID

21100003	3
21100009	3
21100010	1
21100012	1
21100016	3
21100020	3
21100021	3
21100032	1
21100037	1
21100039	3
21100044	1
21100045	1
21100049	3
21100050	1
21100058	1
... (13023 Zeilen unterdrückt)	25795
21137960	1
21137961	1
21137963	3
21137964	3
21137967	2
21137971	1
21137972	1
21137973	1
21137976	1
21137978	2
21137979	3
21137985	2
21137987	1
21137991	3
21138000	1

Diese Kennziffer gruppiert Individuen in ihre Ursprungshaushalte zu Beginn des Panels. Das bedeutet, dass die CID einer Person zeitkonstant gehalten wird und sie immer mit dem Haushalt verbunden wird, zu dem sie initial gehört hat, selbst wenn sie seitdem den Haushalt gewechselt hat.

3 Befragungskontext

wave – Erhebungswelle

1	[1] Welle 1, Teil 1 (2021/22)	13053
2	[2] Welle 1, Teil 2 (2021/22)	6669

3 [3] Welle 2 (2022/23) 6128

Diese Variable identifiziert die (Teil-)Welle, in der die Datenerhebung stattgefunden hat.

4 Statistische Gewichtungsfaktoren

design – Inverse Ziehungswahrscheinlichkeit

70.1461868286133	2
113.673835754395	1
129.841339111328	8
151.398025512695	12
177.21418762207	4
181.577438354492	34
202.015487670898	3
209.343887329102	1
226.846572875977	152
230.22819519043	10
230.803298950195	3
255.753479003906	8
287.660064697266	4
302.295227050781	819
322.924377441406	5
328.682861328125	11
383.379913330078	41
403.530364990234	45
453.192596435547	5060
459.955841064453	155
537.873657226562	91
574.819702148438	824
766.259521484375	1986
806.560424804688	651
905.884948730469	2003
1149.13916015625	9734
1612.62072753906	400
2297.7783203125	3783

Diese Variable enthält die inversen Ziehungswahrscheinlichkeiten (Design-Gewichte) für die SCP-Stichprobe. Das Design-Gewicht berücksichtigt die aus dem Stichprobenziehungsdesign resultierenden ungleichen Ziehungswahrscheinlichkeiten.

Das SCP hat ein zweistufiges Stichprobenziehungsverfahren. Auf der ersten Stufe werden, stratifiziert nach Region und Urbanitätsgrad, Gemeinden gezogen (primary sampling units; PSUs). Auf der zweiten Stufe werden Personen (secondary sampling units; SSUs) aus den Registern dieser Gemeinden gezogen. Generell erfolgte die Stichprobenziehung proportional zur Gemeindegröße. Eine Ausnahme ist die beabsichtigte überproportionale Ziehung in Ostdeutschland. Alle in die Stichprobe gezogenen Personen, die an der Befragung teilnehmen, wurden gebeten ihre weiteren volljährigen Haushaltsmitglieder anzugeben. Diese weiteren Haushaltsmitglieder werden dann ebenfalls befragt. Das führt zu einer höheren

Ziehungswahrscheinlichkeit für größere Haushalte, was ebenso durch die Design-Gewichte berücksichtigt wird.

Aufgrund der Rundung von Nachkommastellen kann es im Codebuch zu einer Zusammenfassung von Werten kommen.

hhf – Hochrechnungsfaktor

89.9704132080078	1
102.496826171875	1
103.33854675293	1
114.508514404297	1
117.915008544922	1
121.074142456055	1
121.479843139648	1
125.444351196289	1
132.868362426758	1
137.444686889648	1
138.32536315918	1
149.049224853516	1
150.570892333984	1
156.400650024414	1
156.875350952148	1
... (25811 Zeilen unterdrückt)	25820
37063.0390625	1
37263.2734375	1
38207.4375	1
38462.2421875	1
38637.3496468707	1
38865.48046875	1
39073.11328125	1
39263.203125	1
39357.0354563788	1
39376.40234375	1
39485.2041726044	1
41794.859375	1
41846.23046875	1
42114.9428418834	1
42127.91796875	1

Diese Variable enthält die Haushalts-Nonresponse-Gewichte für das SCP, die zur Reduzierung von Verzerrungen durch Unit-Nonresponse dienen. Dieser Gewichtungsfaktor ist eine Kombination aus inverser Stichprobenziehungswahrscheinlichkeit, einem Nonresponse-Adjustierungsfaktor und einer Extrapolation zur Zielpopulation der Befragung.

Die inverse Stichprobenziehungswahrscheinlichkeit (siehe DESIGN-Variable) korrigiert für die ungleichen Ziehungswahrscheinlichkeiten in die Panel-Stichprobe (z.B. durch das beabsichtigte Über-Ziehen von Ostdeutschen).

Der Nonresponse-Adjustierungsfaktor korrigiert für Unit-Nonresponse. Dafür wurden Teilnahmewahrscheinlichkeiten basierend auf einem logistischen Regressionsmodell geschätzt. Zur Schätzung der Teilnahmewahrscheinlichkeit wurden Daten aus dem Sampling-Frame

(Altersgruppen, Geschlecht, deutsche Staatsbürgerschaft, Bundesländer) und mikrogeographische Daten als Prädiktoren berücksichtigt. Fehlende Werte in den Prädiktoren wurden mittels multipler Imputation vervollständigt. Prädiktoren wurden durch eine Mischung aus Rückwärtselimination und Vorwärtsauswahl unter Verwendung des Kreuzvalidierungs-Mean-Squared-Errors als Selektionskriterium ausgewählt.

Die Extrapolation basiert auf iterativem proportionalem Fitting (auch als “Raking” bezeichnet) mittels Mikrozensus-Daten zur demographischen Zusammensetzung (Bundesland, Gemeindegröße, Haushaltsgröße) der deutschen Bevölkerung.

Die Gewichte für Erhebungswellen ab Welle 1 Teil 2 wurden durch Multiplikation des initialen Haushalts-Nonresponse-Gewichts der Rekrutierung mit der inversen Teilnahmewahrscheinlichkeit der Ankerperson (AP) an der entsprechenden Erhebungswelle generiert. In Welle 1 Teil 1 wurde diese mittels logistischer Regression geschätzt. Die Prädiktorvariablen umfassen hier zusätzlich Umfragedaten der Vorwellen sowie Interaktionsterme für alle Variablen mit dem Befragentyp (AP vs. Haushaltsmitglieder). (Dieses Modell zur Schätzung der Bleibewahrscheinlichkeit ist dasselbe wie das Modell, das in phrf zur Schätzung der Bleibewahrscheinlichkeit einzelner Befragter verwendet wurde.) In den folgenden Wellen (ab Welle 2) wurde die Teilnahmewahrscheinlichkeit der AP multiplikativ auf der Grundlage von zwei Modellen geschätzt: (1) der Wahrscheinlichkeit, dass ein Haushalt noch zu dieser Umfragewelle eingeladen wird, und (b) der Wahrscheinlichkeit, dass die AP bei Einladung teilnimmt. Wie in Vorwellen wurde Multiple Imputation angewandt, um fehlende Daten zu vervollständigen, und relevante Prädiktorvariablen wurden via Rückwärtselimination und Vorwärtsauswahl ausgewählt. Anschließend wurden die Gewichte unter Verwendung von Mikrozensus-Informationen erneut geraket.

Aufgrund der Rundung von Nachkommastellen kann es im Codebuch zu einer Zusammenfassung von Werten kommen.

5 Inverse Bleibewahrscheinlichkeit

hbleib – Inverse Bleibewahrscheinlichkeit

0	13589
1.04099309444427	1
1.04723525047302	1
1.0492650270462	1
1.04978251457214	3
1.05038011074066	1
1.05090939998627	1
1.05113232135773	1
1.05206882953644	1
1.05332958698273	1
1.05417013168335	1
1.05545318126678	1
1.0565847158432	1
1.05744814872742	1
1.05799996852875	1
... (9884 Zeilen unterdrückt)	12230
5.29517936706543	1
5.31301116943359	1
5.42971420288086	1

5.43652248382568	1
5.45314073562622	1
5.51391649246216	1
5.60983037948608	1
5.65867233276367	1
5.72668218612671	1
5.77419185638428	1
5.92092227935791	1
6.20641374588013	1
6.32693386077881	1
6.40488052368164	1
7.49871492385864	1

Diese Variable enthält die inverse Bleibewahrscheinlichkeit des Haushalts in den Wellen nach der Rekrutierung entsprechend einer Modellierung der Bleibewahrscheinlichkeit der Ankerperson (AP) mittels logistischer Regression. Die Prädiktorvariablen umfassen Erhebungsdaten aus früheren Wellen, einschließlich Interaktionsterme mit dem Befragungstyp (AP vs. Haushaltsmitglied). Fehlende Werte in den Prädiktoren wurden mittels multipler Imputation vervollständigt. Prädiktoren wurden durch eine Mischung aus Backward- und Forward-Selection unter Verwendung des Kreuzvalidierungs-Mean-Squared-Errors als Selektionskriterium ausgewählt.