



Research Institute
Social Cohesion

RDC

SOEP

SCP Documentation

German Social Cohesion Panel

SCP 2021/22 W1 Codebook PHRF: Weights for Persons (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., Richter, D., Teichler, N., Traummüller, R., & Zinn, S. (2024). SCP 2021/22 W1 - Codebook PHRF: Weights for Persons (German). *German Social Cohesion Panel 2021/22 - Wave 1*. Bremen and Berlin: RDC-RISC/SOCIUM, SOEP/DIW Berlin. doi:10.60532/scp.2021_22.w1.v1

- ▶ **Authors:** Olaf Groh-Samberg, Julian B. Axenfeld, Jean-Yves Gerlitz, Carina Cornesse, Martin Kroh, Holger Lengfeld, Stefan Liebig, Lara Minkus, Jost Reinecke, David Richter, Nils Teichler, Richard Traummüller, Sabine Zinn
- ▶ **Contributors:** Cosima Adams, Anton Bochert, Martin Gerike, Josefine Kuhrmeier, Anna-Tabea Müller, Eric Nissen, 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.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 Codebook PHRF: Weights for Persons (German)

Inhaltsverzeichnis

1	Allgemeine Informationen	2
2	Identifikatoren	2
	pid – Unveränderliche Personen-ID	2
	hid – Aktuelle Haushalts-ID	3
	cid – Ursprüngliche Haushalts-ID	3
3	Befragungskontext	4
	wave – Erhebungswelle	4
4	Statistische Gewichtungsfaktoren	4
	design_ap – Inverse Ziehungswahrscheinlichkeit Ankerpersonen	4
	design – Inverse Ziehungswahrscheinlichkeit	5
	phrf – Hochrechnungsfaktor	6
5	Inverse Bleibewahrscheinlichkeit	8
	pbleib – Inverse Bleibewahrscheinlichkeit	8

1 Allgemeine Informationen

Der PHRF-Datensatz enthält Gewichtungsfaktoren für die einzelnen Befragten im SCP. Jede Person (PID), die in einer bestimmten Erhebungswelle (WAVE) an der Befragung teilgenommen hat, hat eine Zeile im Datensatz.

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

pid – Unveränderliche Personen-ID

2110000301	2
2110000302	1
2110000901	2
2110000902	2
2110001001	1
2110001201	1
2110001601	2
2110002001	2
2110002101	2
2110003201	1
2110003701	1
2110003901	2
2110004401	1
2110004402	1
2110004403	1
... (16997 Zeilen unterdrückt)	26150
2113796701	1
2113796702	2
2113796703	2
2113797101	1
2113797102	2
2113797201	1
2113797301	1
2113797601	1
2113797801	2
2113797901	2
2113798501	2
2113798701	1
2113799101	2
2113799102	2
2113800001	1

Die PID ist die unveränderliche Kennziffer einer Person, die über die Zeit und über alle Datensätze identisch gehalten wird.

hid - Aktuelle Haushalts-ID

21100003	3
21100009	4
21100010	1
21100012	1
21100016	2
21100020	2
21100021	2
21100032	1
21100037	1
21100039	2
21100044	3
21100045	1
21100049	2
21100050	1
21100058	1
... (13023 Zeilen unterdrückt)	26135
21137960	1
21137961	1
21137963	4
21137964	4
21137967	5
21137971	3
21137972	1
21137973	1
21137976	1
21137978	2
21137979	2
21137985	2
21137987	1
21137991	4
21138000	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 verändern kann, zum Beispiel wenn ein erwachsenes Kind aus dem elterlichen Haushalt auszieht und einen eigenen Haushalt eröffnet.

cid - Ursprüngliche Haushalts-ID

21100003	3
21100009	4
21100010	1
21100012	1
21100016	2
21100020	2
21100021	2

21100032	1
21100037	1
21100039	2
21100044	3
21100045	1
21100049	2
21100050	1
21100058	1
... (13023 Zeilen unterdrückt)	26135
21137960	1
21137961	1
21137963	4
21137964	4
21137967	5
21137971	3
21137972	1
21137973	1
21137976	1
21137978	2
21137979	2
21137985	2
21137987	1
21137991	4
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)	17027
2	[2] Welle 1, Teil 2 (2021/22)	9168

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

4 Statistische Gewichtungsfaktoren

design_ap – Inverse Ziehungswahrscheinlichkeit Ankerpersonen

0	6623
905.88493688161	4547
905.884936881611	1562
1612.62063492063	917

2297.77831821929 12546

Diese Variable enthält die inversen Ziehungswahrscheinlichkeiten (Design-Gewichte) für die initiale Stichprobe von Ankerpersonen. Das Design-Gewicht berücksichtigt die aus dem Stichprobenziehungsdesign resultierenden ungleichen Ziehungswahrscheinlichkeiten, die aus der überproportionalen Ziehung von Personen in Ostdeutschland resultieren. Diese Design-Gewichte sollen für Analysen verwendet werden, die nur die initiale Ankerpersonenstichprobe ohne die weiteren Haushaltsmitglieder nutzen.

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.

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

design – Inverse Ziehungswahrscheinlichkeit

70.1461868286133	2
113.673835754395	1
129.841339111328	19
151.398025512695	23
177.21418762207	4
181.577438354492	47
202.015487670898	16
209.343887329102	1
226.846572875977	200
230.22819519043	28
230.803298950195	2
255.753479003906	8
287.660064697266	3
302.295227050781	964
322.924377441406	5
328.682861328125	11
383.379913330078	57
403.530364990234	53
453.192596435547	5246
459.955841064453	217
537.873657226562	108
574.819702148438	1074
766.259521484375	2460
806.560424804688	682
905.884948730469	1525
1149.13916015625	10200
1612.62072753906	311
2297.7783203125	2928

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.

phrf - Hochrechnungsfaktor

138.076012832869	1
180.275392341203	1
191.814935987255	1
198.18234981009	1
202.519416418032	1
205.922105072994	1
225.538315268451	1
225.66939124434	1
230.998437933074	1
231.575753584113	1
242.620946761035	1
250.467474531022	1
252.662663340351	1
254.528631438953	1
255.406094077172	1
... (26124 Zeilen unterdrückt)	26165
60540.2697879398	1
60553.0048254172	1
60651.5267019675	1
60926.7044585749	1
61361.9541063377	1
61490.091541196	1
61693.9423795728	1
62316.1498544427	1
63344.4531750494	1
63985.7215170963	1
64328.3586230575	1
64630.1889132602	1
65160.1943149763	1
66936.6426181433	1
67220.5262768564	1

Diese Variable enthält die individuellen 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. Für seine Berechnung wurden Teilnahmewahrscheinlichkeiten durch eine Verkettung mehrerer Regressionsmodelle geschätzt:

1. Ein logistisches Regressionsmodell der Teilnahmewahrscheinlichkeit der Ankerperson (AP), in der Sampling-Frame-Daten und zusätzliche mikrogeographische Daten als Prädiktoren berücksichtigt wurden. (Dieses Modell ist dasselbe wie das, das zur Generierung von hhrf genutzt wurde.) 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.
2. Ein fraktionelles Regressionsmodell des Anteils der Haushaltsmitglieder (HM), die durch die AP zur Teilnahme an der Studie genannt wurden. Dadurch wird eine Untererfassung der HM durch die zugehörige AP berücksichtigt. Hierbei werden zusätzlich zu Sampling-Frame- und mikrogeographischen Daten auch Befragungsdaten der AP als Prädiktoren genutzt. Wie in Modell (1) wurden fehlende Werte durch multiple Imputation ergänzt und relevante Prädiktorvariablen wurden via Rückwärtselimination und Vorwärtsauswahl ausgewählt.
3. Ein logistisches Regressionsmodell der Teilnahmewahrscheinlichkeit des HM. Hierbei wurden zusätzlich zu Sampling-Frame-Daten, mikrogeographischen Daten sowie Befragungsdaten der Ankerperson auch das Alter des HM, das Geschlecht des HM sowie die Beziehung des HM zur AP als Prädiktoren genutzt. Wie in Modell (1) wurden fehlende Werte durch multiple Imputation ergänzt und relevante Prädiktorvariablen wurden via Rückwärtselimination und Vorwärtsauswahl ausgewählt.

Für AP lässt sich die Teilnahmewahrscheinlichkeit direkt aus Modell (1) ableiten, während sich die Teilnahmewahrscheinlichkeit für HM aus der Multiplikation der vorhergesagten Wahrscheinlichkeiten der Modelle (1) bis (3) zusammensetzt.

Die Extrapolation basiert auf iterativem proportionalem Fitting (auch als "Raking" bezeichnet) mittels Mikrozensus-Daten zur demographischen Zusammensetzung (Alter, Geschlecht, deutsche Staatsbürgerschaft, Ost- vs. Westdeutschland) der deutschen Bevölkerung.

Die Gewichte für Erhebungswellen ab Welle 1 Teil 2 wurden durch Multiplikation des initialen Nonresponse-Gewichts der Rekrutierung mit der inversen Teilnahmewahrscheinlichkeit an der entsprechenden Erhebungswelle, geschätzt mittels logistischer Regression, generiert. Die Prädiktorvariablen umfassen hier zusätzlich Umfragedaten der Vorwellen, einschließlich Interaktionsterme für alle Variablen mit dem Befragentyp (AP vs. HM). (Dieses Modell zur Schätzung der Bleibewahrscheinlichkeit ist dasselbe wie das Modell, das in hhrf zur Schätzung der Bleibewahrscheinlichkeit der Haushalte verwendet wurde.) Auch hier 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

pbleib – Inverse Bleibewahrscheinlichkeit

0	17027
1.11346220970154	1
1.12490344047546	1
1.12930381298065	1
1.13066411018372	1
1.13112485408783	1
1.13122057914734	1
1.13646674156189	1
1.13921761512756	1
1.1408588886261	1
1.14317750930786	1
1.14333009719849	1
1.14349722862244	1
1.14467191696167	1
1.14593815803528	1
... (8960 Zeilen unterdrückt)	9139
6.88516664505005	1
6.88868188858032	1
6.91688442230225	1
6.92782402038574	1
7.00820970535278	1
7.04780006408691	1
7.04848480224609	1
7.21949863433838	1
7.34550809860229	1
7.46248769760132	1
7.49639701843262	1
7.84853172302246	1
8.08865928649902	1
9.17348957061768	1
10.321307182312	1

Diese Variable enthält die individuelle inverse Bleibewahrscheinlichkeit in den Wellen nach der Rekrutierung entsprechend einer Modellierung durch logistische Regression. Die Prädiktorvariablen umfassen Erhebungsdaten aus früheren Wellen, einschließlich Interaktionsterme mit dem Befragungstyp (Ankerperson 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.