

# Survey Data Analysis Final Report

*Gerbrich Ferdinands*

*Sanne Meijering*

*Hanne Oberman*

## Abstract

In this technical report, a single wave of the multi-year panel study ‘Understanding Society’ is analyzed. More specifically, this research project focussed on the use of design and adjustment weights for estimating population quantities from a subset of the ‘Understanding Society’ data: the ‘Innovation Panel’. The data was analyzed using the R package ‘survey’, and base R operations. **[Add sentences about results and discussion!]**.

*Key terms:* survey package, design weights, adjustment weights, nonresponse.

## Table of Contents

Introduction

Methods

Results

Discussion

References

Appendix I: Research Proposal

Appendix II: R Script

## Introduction

To estimate population quantities from sample data, broadly two strategies can be employed: model-based inference and design-based inference (Lohr, 2010). In this technical report, the effect of the latter on population estimates is investigated, by means of analyzing a single wave of Great Britain’s multi-year panel study ‘Understanding Society’. This annual household panel survey is used to assess topics like housing, education, health, and employment among the general population of the United Kingdom. More specifically, we focus on the a subset of first wave, conducted between 2008 and 2009: the ‘Innovation Panel’ (IP).

The purpose of the IP was to enable methodological research, before initiating the main ‘Understanding Society’ survey. As a forerunner, the target sample size of the IP was 1500 households across Great Britain, whereas the intended sample size for the main survey was 40,000 households within the United Kingdom (Boreham, & Constantine, 2008). Therefore, the sampling frame of the IP, The Postcode Address File, encompassing England, Wales and Scotland. Other regions of the United Kingdom like Northern Ireland were not part of the target population.

To obtain the sample from the sampling frame, multistage stratified cluster sampling was used. Postal sectors acted as primary sampling units (PSUs) in the cluster design, and were sorted by Government Office Region. Within these regions, strata were constructed based on a.o. non-manual occupations, ethnic minority density, and population density. After stratification, a sample was drawn from the sampling frame (the Postcode Address File), with probabilities proportional to the number of postal delivery points within each postcode sector (the stratified PSUs). Per postcode sector, 23 addresses were selected, yielding a total of 2760 addresses. However, interviews were only achieved in 1489 households (Boreham, & Constantine, 2008).

**[Add paragraph about Non-sampling, sampling, and nonresponse errors].**

The aim of this research project is to investigate several properties of the IP data, with respect to the sampling procedure, the use of weights in survey data analysis, and to estimate population quantities from the sample.

The report is based on eight research questions (RQs) that were provided in the research proposal (Appendix I: Research Proposal). The RQs are structured according to three topics: sampling, weighting, and estimation. In the sampling part, we investigate how within-household selection of respondents was performed (RQ1), and how many people within the sample personally conducted an interview (RQ6). The part concerning survey weights encompasses a theoretical explanation for the variation between design weights across observations (RQ2), computation of enumerated weights (RQ3), and an effort to suggest variables that could be used to construct nonresponse weights (RQ8). Finally, in the third part population quantities are estimated: the relationship between age and household composition (RQ5), and the proportion of employed people of working age in the population with and without accounting for nonresponse (RQ4 and RQ7, respectively). The exact formulation of the research questions can be found in Appendix I: Research Proposal.

## Methods

The data in this research project was obtained by ‘The National Centre for Social Research’ from the United Kingdom (Boreham, & Constantine, 2008, p. 3), and retrieved via Dr. Peter Lugtig (Utrecht University). Data analysis is performed using the statistics freeware R (R Core Team, 2003). Preprocessing of the data included recoding the variable ‘a\_dvage’ as numeric instead of factor, and removal of irrelevant variables from the dataframe (see Appendix II: R Script).

To obtain accurate estimates from survey data with a multistage sampling design, regular SRS estimators do not suffice (Lohr, 2010). Therefore, the R package ‘survey’ (Lumley, 2018) is used. By specifying the survey design and appropriate weights, estimates of population quantities and their variances are adjusted for the complex nature of the survey data. The survey package functions ‘svydesign’, and ‘svymean’ are employed to estimate population quantities. Moreover, stepwise logistic regression is performed to investigate which variables might be used to construct nonresponse weights. This exploratory analysis is conducted using the R package ‘MASS’ (Venables, & Ripley, 2002).

## Results

### Part I: Sampling

#### Within-household selection of respondents (RQ1)

For the Understanding Society study, up to three households per dwelling and up to three people per household were selected for interviews. If more than three households were found at a single dwelling, or if a single household contained more than three people, the Kish Grid method was used to select households or people at random.

This method considered to lead to a random sample of household members, and avoids selection bias in the survey. When using the Kish Grid method, the researcher has to list all eligible household members and assign each of them a number  $i$ . The Kish Grid then helps you pick a random person  $i$  within the  $n^{\text{th}}$  household you’re investigating, see Figure 1.

It is a popular method because when carried out correctly, it leads to (almost) equal probability sampling, something other selection methods do not obtain. This is because of the way the grid is set up, each household member has an equal chance of being selected. Moreover, the only information you need in order to select respondents is a list containing the names of the persons in the household you want to sample, because selection is based on the number you’ve assigned to them. Other selection methods usually need more information, like the date of birth of the household members.

#### Proportion of personal interviews within the sample (RQ6)

Within the data, there is a proportion of sample members that did not personally take part in an interview, potentially leading to nonresponse bias. However, not all of these observations count as nonresponse. Of the 3600 individuals in the sample, there are 459 children under the age of ten. These sample members are

	Eligible People							
Household	1	2	3	4	5	6	7	8+
1st	1	1	1	1	1	1	1	1
2nd	1	2	2	2	2	2	2	2
3rd	1	1	3	3	3	3	3	3
4th	1	2	1	4	4	4	4	4
5th	1	1	2	1	5	5	5	5
6th	1	2	3	2	1	6	6	6
7th	1	1	1	3	2	1	7	7
8th	1	2	2	4	3	2	1	8
9th	1	1	3	1	4	3	2	1
10th	1	2	1	2	5	4	3	2

Figure 1: Kish Grid. *Note.* Reprinted from <https://www.statisticshowto.datasciencecentral.com/wp-content/uploads/2017/09/kish-grid.png>

enumerated individuals, but ineligible for an interview. Accordingly, the number of eligible individuals in the sample is 3141. Of these 3141 eligible individuals, only 2656 respondents personally completed an interview. This number consists of 2399 adults who performed a full interview, and 257 interviews with children. Taken together, 485 cases of (partial) unit nonresponse exist within the subset of eligible individuals.

Moreover, there are 169 sample members for whom a proxy interview was conducted with a close relative or cohabitant (generally a spouse, partner or adult child; Understanding Society, ‘Questions about variables’, N.D.). Questions answered on the nonresponders’ behalf were used to estimate factual information about the nonrespondent. Information about their beliefs is not collected by proxy, and are therefore indicated as missing (item nonresponse).

Because (partial) nonresponse can lead to nonresponse bias, it is important to investigate the effect of nonresponse on the estimates generated above. Ideally, we would both compare respondents to partial nonresponse, and unit nonresponse. Because of time constraints within this research project, we only investigate unit nonresponse. Of the 3600 individuals in the sample, there are 459 children under the age of ten. These sample members are enumerated individuals, but ineligible for an interview. Therefore, the number of eligible individuals in the sample is 3141. Of these 3141 eligible individuals, only 2656 respondents personally completed an interview. This number consists of 2399 adults who performed a full interview, and 257 interviews with children. Taken together there are 485 cases of (partial) unit nonresponse within the subset of eligible individuals.

Moreover, there are 169 sample members for whom a proxy interview was conducted with a close relative or cohabitant (generally a spouse, partner or adult child; Understanding Society, ‘Questions about variables’, N.D.). Questions answered on the nonresponders’ behalf were used to estimate factual information about the nonrespondent. Information about their beliefs is not collected by proxy, and are therefore indicated as missing (item nonresponse).

Because (partial) nonresponse can lead to nonresponse bias, it is important to investigate the effect of nonresponse on the estimates generated above. Ideally, we would both compare respondents to partial nonresponse, and unit nonresponse. Because of time constraints within this research project, we only investigate how weights may be constructed to adjust for unit nonresponse (see Part III: Estimation).

## Part II: Weighting

### Variation in design weights (RQ2)

Following from the sampling design, all sample members were assigned a design weight. As Table 1 shows, the variance in the design weight variable is 0.02. With a mean of 1.01, the coefficient of variation is only 0.14.

Table 1. *Descriptive statistics of design weight variable ‘a\_psnenip\_xd’.*

Mean	SD	Min.	Max.
1.01	0.14	1	4

The low variance within the design weight variable can be due to rescaling of the weights, as explained in the IP User Guide (p. 56): “Each set of weights has been scaled by a constant factor to produce a mean of one amongst cases eligible to receive the weight”. That is, the design weights are standardized on one. This procedure does, however, not affect the coefficient of variation.

Moreover, the maximum value of this variable suggests trimming of weights greater than 4. Trimming, truncating, or smoothing of weights are techniques to inhibit the effect of a single observation on population estimates and its variance. Trimming may yield biased estimates, but is shown to decrease the mean square error (MSE; Lohr, 2010). It is therefore reasonable to assume that weight trimming was performed within the current dataset.

A histogram of the design weight variable (see Figure 2) shows a different reason for the low variance: all but 26 respondents have a design weight of exactly 1. This can be attributed to the sampling design, in which, among other things, population density within PSUs was used to construct strata. With that, inclusion probabilities of households were roughly equal within strata. The resulting design weights thus represent the discrepancy between the number of observation units in the population, and the proportion of households per stratum selected from the sampling frame.

### Computation of enumerated weights (RQ3)

Apart from design weights, each respondent is also assigned an enumerated weight (the variable ‘a\_psnenip\_xw’). This post-stratified weight was calculated using the design weights, sex, age and government office region. It is known that age was split up in seven groups and the government office regions were split up in four, but which ones these were is unknown.

In order to find out, we first turned age into a numeric variable and created dummy variables for all government office regions, with Scotland being the baseline. Sorting the coefficients from lowest to highest yielded the following graph:

It is not completely clear which values should be grouped together into four groups, but there are two groups of three regions and one group of two with near-identical weights. As can be seen on the map below, where the regions are numbered from the lowest to the highest weight, the grouping is rather strange.

The baseline and the other two areas with a weight close to zero are Scotland, London and the West Midlands. Scotland and London are prime candidates to be categories of their own, as Scotland has a strong regional identity and London is the capital of the country. From a data-driven perspective, combining the West Midlands with any region other than London or Scotland does not make sense, as the weight difference between those region and all other regions is relatively big. From a theory-driven perspective, we cannot think of any good reason to combine the West Midlands with London or Scotland: it is not filled with cities nor is it close enough to Scotland to assume a similar population.

When assuming that the regions with a near-identical weight are in the same group and Scotland and London are two separate categories, only a few possible groupings remain:

1. South East, South West, and North East England are the third category, with Wales and the rest of England being the last category.

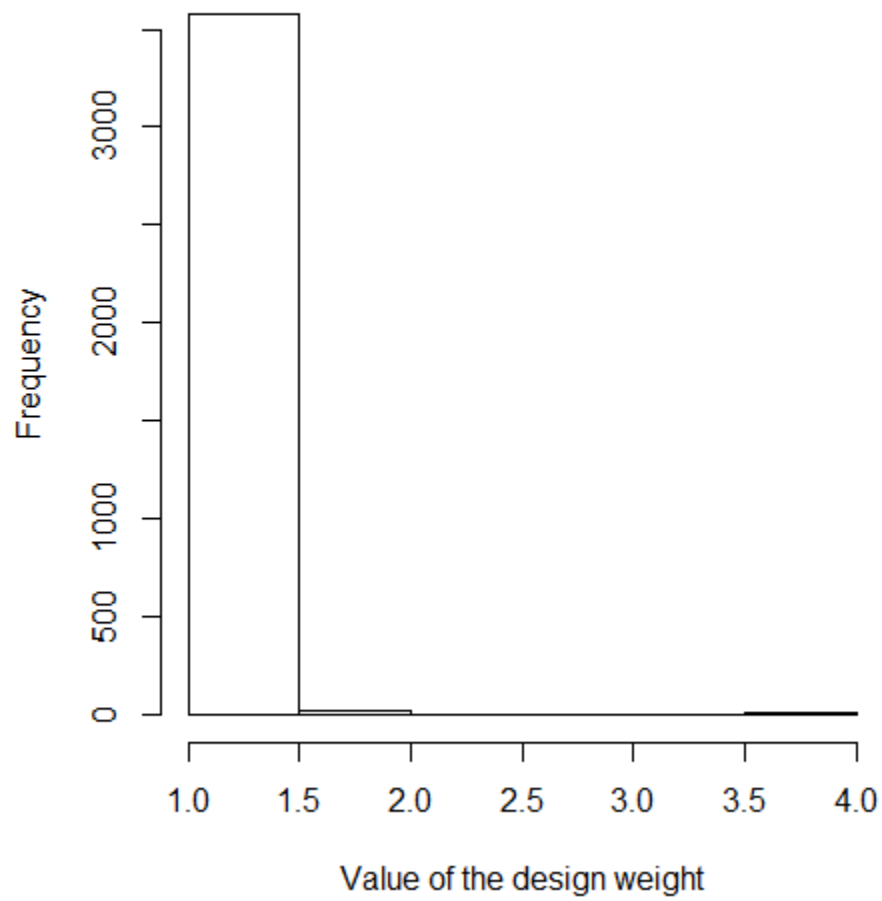


Figure 2: Histogram of design weight variable 'a\_psnenip\_xd'.

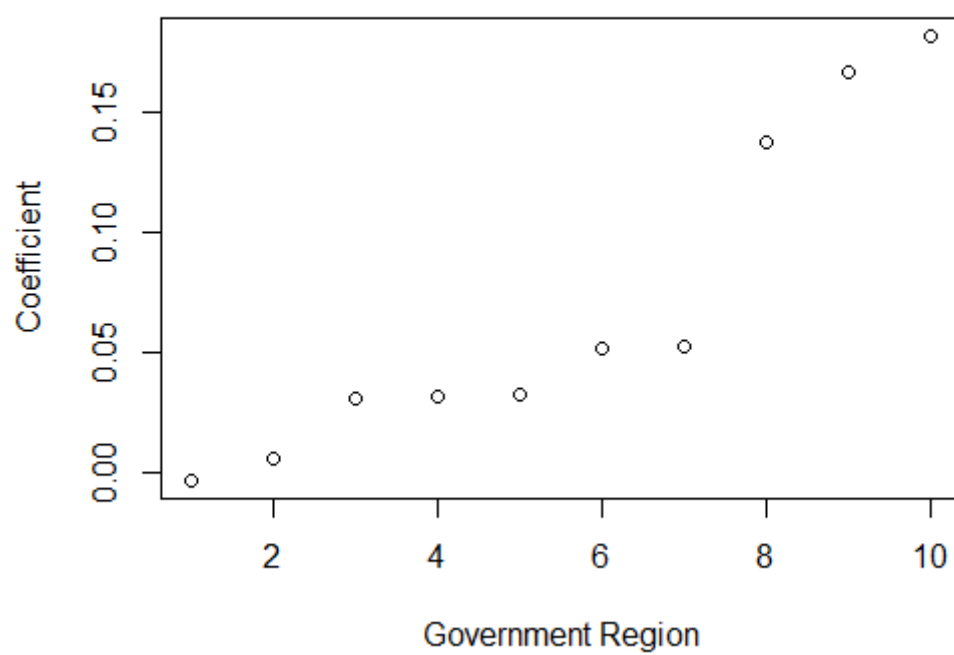


Figure 3: Coefficients per Government Region.



Figure 4: Government Regions in Great Britain.

This makes some sense theoretically: the three are the very north and south of England, with the remainder being in the middle. Weight-wise however, combining the West Midlands with the areas with far lower weights does not make sense.

2. The same as above, but with the West Midlands grouped with London.

This is more logical weight-wise than the first option, but from a theoretical perspective it is strange to combine the city with a rural area.

3. The North West, West Midlands and East of England are the third category, with Wales and the rest of England being the last category.

Weight-wise it is just as logical or illogical as the first option, but from a theoretical standpoint taking three areas in the middle seems to have less of a basis than taking the outer areas like is done in the first option.

4. The West Midlands alone being the third category, with the rest of England being the fourth.

This fits weight-wise, but why the West Midlands should be taken separately is a mystery. It is however less absurd than combining it with London: the West Midlands may have unique features that we do not know about, but we do know that it is not a capital like London.

In conclusion, the fourth option seems the most likely, as the grouping is logical from a data-driven perspective and while it does not quite make sense theoretically, there is no obvious reason to assume that it is incorrect. The best alternative is the first option, as it makes the most sense from a theoretical perspective.

Considering age groups, the most logical categories are 0-20, 21-30, 31-40, 41-50, 51-70, 71-80, 81+. Theoretically speaking, the 51-70 category is a bit strange, but due to the rather convincing difference in weights between the 71-80 group and the neighbouring groups, two of the categories between 21 and 70 must be joined, and 51-70 is the most uniform weight-wise.

A regression analysis revealed which groups were under- and overrepresented in the sample (compared to the population of Great Britain), with higher weights being assigned to underrepresented groups. While the weights are mostly decided by the design weights, the post-stratification does alter them.

Females have slightly lower weights than males, indicating that they were overrepresented. For government office region groups, the people from England (excluding London and the Western Midlands) and Wales are the baseline and have the lowest of all weight, and thus were also overrepresented, with Scotland being the most underrepresented, followed by London and then the Western Midlands. With regards to age, the 0-20 group was surprisingly the most overrepresented, forming the baseline for this variable. The 21-30 and 31-40 groups were the most underrepresented and thus had the highest weights, followed by the 80+ group, the 51-70 group, the 41-50 group and finally the 71-80 group, which was by far the most underrepresented.

## **Construction of nonresponse weights (RQ8)**

Nonresponse can seriously threaten the quality of estimates if the nonresponders differ from the responders in the study. We can adjust for this by the use of auxiliary variables that hold information about the entire sample. We can compare the distribution of the auxiliary variables of the population with its response distribution, and assess whether the response is representative for the whole population.

If the distributions differ, we must conclude that nonresponse has led to a nonrepresentative sample. We can adjust for this by constructing nonresponse weights from the auxiliary variables. Then the weighted sample is representative with respect to the auxiliary variables used. Different choices of auxiliary variables lead to different weights. In order for the weights to be useful, it is important to use powerful auxiliary variables (Brick, 2013).

The IP User Guide states that “all models used to predict response propensities as described in the Technical Details are fitted using stepwise backward logistic regression with  $p=0.05$ ” (p. 56). Compared to other selection strategies, this is a fairly simple one. Considering these other strategies for selecting auxiliary variables for weighting adjustments for nonresponse, we advise to look into Schouten’s (2007) proposed strategy of a forward inclusion-backward elimination algorithm and/or Sørensen and Lundström’s (2010)



paper. We opt for a stepwise backwards logistic regression of because this method is similiar to, yet easier to apply than the selection strategies described in Schouten (2007) S??rndal and Lundstr??m's (2010). Most important, this selection strategy was also used in constructing the original weights in the dataset. With this method, we hope to identify variables that predict nonresponse.

There are two stages of nonresponse we would have wanted to consider: \* person nonresponse: whether a personal interview was conducted yes or no; \* household nonresponse: defined as completion of at least the household grid.

However, as the levels of variable 'a\_hhresp\_dv' show, all households within the sample completed at least the household grid. We do not have any information about the 1271 adresses that were selected from the sampling frame but did not take part in an interview. Therefore, we only account for individual unit nonresponse, which occurred in 13.5% of the observations in the sample.

The backwards logistic regression model shows us that the variables ... are useful in predicting nonresponse. The predictive accuracy of the stepwise model is ..., where for the full model this was ..

## Part III: Estimation

### Relationship between age and household composition (RQ5)

To investigate the relationship between age and household composition, three steps were performed: choosing which age variable would be used, constructing a new variable to represent household composition, and including the appropriate weight in the analysis.

Firstly, we investigated which variable to choose representing the age of respondents. There are four variables in the data expressing the age of respondents: a\_agegr5\_dv, a\_agegr10\_dv, a\_agegr13\_dv, a\_dvage. The first three are categorical variables, divided into groups as seen in the UK Labour Force Survey. The fourth variable, a\_dvage, is continuous.

The first categorical variable, a\_agegr5\_dv, contains 15 categories indicating age in five-year partitions up to '70 years or older'. The variable a\_agegr10\_dv is highly similar, indicating age in ten-year intervals, and consisting of 8 different categories between ages zero and seventy-plus. Finally, a\_agegr13\_dv consists 13 levels, with unequal age intervals under the age of 20 (namely 0-15, 16-17, and 18-19 years old), five-year intervals between the ages 20 and 64, and a sigle interval containing all ages from 65 years old upward.

The continous variable a\_dvage denotes the age of respondents in completed years at the moment the interview took place. This was computed by comparing the date of birth and the interview date, or using an estimate (see estimated age variable a\_ageest (n = 22), and imputed age variable a\_ageif (n = 55)).

To investigate the relationship between age and household composition, it is necessary to know which households consist of children living at home, and which don't. Presumably, the largest variation between children living at home and those living independently occurs within the age group 15-25 year-olds. If we would use the ten-year intervals to indicate age, we might get less precise estimates, because of the large variance wihtin the age category '10-19 years old'. Therefore, we decided to use an age variable with smaller intervals. Because of the unequal intervals in the age variable with 13 categories, it should be analysed as ordinal variable. To interpret the relationship with household composition, this is not ideal. Ultimately, we decided to use either a\_dvage or a\_agegr5\_dv.

Secondly, to construct a new variable denoting household composition, we used Eurostat as reference. Eurostat distinguishes three types of variables used to descripte household composition:

- Household size;
- Number of children;
- Household type, with levels: couple (with/without children), single adult (with/without children), other type of household (with/without children).

All three household composition indicators were constructed from the data and included per household into a household composition matrix (see Appendix II: R Script).

Thirdly, we determined which weighting variable is appropriate to investigate the relationship between age and household composition. We want to use a weight per person (PSN), as it is impossible to calculate age per household. We also want to use the enumeration grid weights (EN), as both age and household composition have no missing values. We want the weights of the first wave (IP), as that is the one we are investigating. Finally, we want to use the stratified weights (XW) rather than the design weights, as we are investigating population values and not sample values. Thus, we want to use the weight `a_psnenip_xw`, as this is the only weight that covers the correct data in the correct way.

The relationship between age and household composition is the following: ADD RESULTS!!!

#### **Proportion of employed people of working age (RQ4)**

43,3% of the population is employed, with a 95% confidence interval of 41.4%-45.2%.

#### **Proportion of employed people of working age, accounting for nonresponse (RQ7)**

After post-stratification (as we also included post-stratification in the previous analysis), 46% of the people who did a personal interview is employed, with a confidence interval of 44.9%-49.3%. For the previous question, where all but two people were included, the estimated proportion of the population that is employed was 43,3%, with a 95% confidence interval of 41.4%-45.2%. As neither mean is contained in the 95% confidence interval of the other, it can be concluded that there is a difference between respondents and non-respondents in employment, with respondents being more likely to be employed than non-respondents.

## **References**

- Boreham, R., Constantine, R. (2008). Understanding Society Innovation Panel Wave 1: Technical Report. National Center for Social Research.
- Eurostat. (2018). Household composition statistics [Explainer]. Retrieved via <https://ec.europa.eu/eurostat/statistics-explained/pdfscache/29071.pdf>.
- Kumar, Rohit. (2014). Respondent Selection Methods in Household Surveys. Jharkhand Journal of Development and Management Studies. XII. 5701-5708.
- Lohr, S. (2010). Sampling: Design and Analysis (second edition). Brooks/Cole, Boston, USA.
- Lumley, T. (2018). Survey: analysis of complex survey samples [R package, version 3.34]. Retrieved via <https://cran.r-project.org/web/packages/survey/>.
- R Core Team (2013). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. Retrieved via <http://www.R-project.org/>.
- Sa?rndal, C.E. and Lunsdström, S. (2010). Design for Estimation: Identifying Auxiliary vectors to reduce nonresponse bias. Survey Methodology, 36, 131-144.
- Schouten, B. (2007). A Selection Strategy for Weighting Variables Under a Not-Missing-at-Random Assumption. Journal of Official Statistics, 23, 51-68.
- Venables, W. N. & Ripley, B. D. (2002) Modern Applied Statistics with S. Fourth Edition. Springer, New York. ISBN 0-387-95457-0

## **Appendix I: Research Proposal**

Understanding Society is a large, multi-year panel study of the population in the United Kingdom and Northern Ireland. The dataset that you will analyse is wave 1 of the Understanding Society Innovation Panel, which is used for testing instruments. The Innovation panel (IP) uses the same sampling procedure as the main Understanding Society Survey, apart from the fact that Northern Ireland is not included. Sampling is done using a multistage stratified cluster sample. For more information, see <https://www.understandingsociety.ac>.

uk/sites/default/files/downloads/documentation/innovation-panel/user-guides/ip\_user\_guide.pdf Within this file, further references are made to the sampling and weighting guidelines. Do read these, because this will be important for the rest of the analyses. The data file you need for this assignment is ‘understanding society innovation panel wave A.RDS’.

Questions for report: (for questions 1-5), please ignore nonresponse errors. 1. The Innovation panel would interview multiple people from every household, but in the case households are very large, a Kish Grid is used. Can you explain why the Kish Grid is a popular method to do within-household selection of respondents?

2. Despite the complex data structure, the design weights as included in the variable ‘a\_psnenip\_xd’ have a very low variance. Can you explain why this is the case?
3. Apart from the design weight, an enumerated design weight ‘a\_psnenip\_xw’ is also included. Please show using syntax in R how you can calculate this weight.
4. What is the proportion of employed people of working age (15-64) in the population? (use variable ‘a\_employ’)
5. Please investigate the relationship between age and household composition for the population of Great Britain. Take the following steps:
  - Investigate which age variable you should use (and whether choosing a particular age variable at all matters)
  - Construct a new variable ‘household\_composition’ from the following four variables: ‘a\_livesp\_dv, a\_cohab\_dv, a\_single\_dv, a\_mastat\_dv’. You may also use information from across household members to do this.
  - Which weight should you use here? From here on, take nonresponse error into account.
6. Not all sample members you included in questions 1-5 participated in the survey. Use the variable ‘a\_ivfio’ to investigate how many people personally conducted an interview.
7. Repeat your analysis under question 4. Use only the people who personally did an interview. What is your conclusion? Is there a difference between respondents and nonrespondents when it comes to their employment?
8. Using only variables from this dataset, can you find variables which could potentially be useful for constructing nonresponse weights? Why these variables?