

NIASRA

NATIONAL INSTITUTE FOR APPLIED
STATISTICS RESEARCH AUSTRALIA



***National Institute for Applied Statistics Research
Australia***

University of Wollongong, Australia

Working Paper

10-18

**How to use Replicate Weights in Health Survey Analysis
using the National Nutrition and Physical Activity Survey as
an Example**

Carole L. Birrell, David G. Steel, Marijka J. Batterham and Ankur Arya

*Copyright © 2018 by the National Institute for Applied Statistics Research Australia, UOW.
Work in progress, no part of this paper may be reproduced without permission from the Institute.*

National Institute for Applied Statistics Research Australia, University of Wollongong,
Wollongong NSW 2522, Australia Phone +61 2 4221 5435, Fax +61 2 4221 4845.
Email: karink@uow.edu.au

Abstract

Objective

To conduct nutrition related analyses on large-scale health surveys, two aspects of the survey must be incorporated into the analysis: the sampling weights and the sampling design; a practice which is not always observed. The aim of this paper is to compare three analyses: 1. unweighted; 2. weighted but not accounting for the complex design; and 3. weighted and accounting for the complex design.

Design

Analyses with and without the proper use of sampling weights and replicate weights are conducted using Stata. Descriptive statistics are computed and a logistic regression investigating whether chosen explanatory variables are associated with being overweight/obese.

Setting

Cross-sectional health survey with complex sampling design when replicate weights are supplied.

Subjects

Responding adults from the National Nutrition and Physical Activity Survey (NNPAS) part of the Australian Health Survey (2011-13).

Results

An unweighted analysis gives both a biased estimate and incorrect standard errors. Adjusting for the sampling weights give unbiased estimates but incorrect standard errors. Incorporating both the sampling weight and design results in unbiased estimates and the correct standard error. This can have effects on interpretation, for example, the odds ratio for being a current smoker in the unweighted analysis was 1.20 (95%CI 1.06,1.37) $t=2.89$, $P=0.004$, suggesting a statistically significant relationship with being overweight/obese. When the sampling weights and complex design are incorporated the results are no longer significant, odds ratio =1.06 (95% CI 0.89,1.27) $t=0.71$, $P=0.480$.

Conclusions

Correct incorporation of the sampling weights and design are crucial for valid inference from survey data.

33 **Keywords:** Complex survey design, replicate weights, survey sampling, sampling weights,
34 health surveys, BMI.

35 1 Introduction

36 Many nutrition and public health researchers make use of data obtained from
37 large-scale surveys to estimate relationships between nutrition and the health status of
38 the population and particular sub-groups, and to inform health policies. The
39 Australian Health Survey (AHS) and the United States National Health and Nutrition
40 Examination Survey (NHANES) are two health surveys collected using a complex
41 survey design. Complex sampling procedures may involve use of many design features
42 such as geographic stratification, multistage sampling involving clustering, and the
43 disproportionate sampling of certain ethnic or age groups. In order to validly
44 generalise the results to the relevant population, the study design features must be
45 incorporated into the estimation and analysis.

46 Analysis of data resulting from a complex sample survey to produce unbiased
47 estimates and estimated standard errors (SEs) which account for the sampling design
48 can be complicated⁽¹⁾. It requires the use of the individual's sampling weight and the
49 sampling design variables; the resulting estimates are called *design-based* estimates⁽²⁾.
50 The sampling weight is based on the inverse of the probability of selection and will
51 often vary considerably between individuals, due to the sampling design and
52 post-survey adjustments. It can be considered as the number of units (such as
53 individuals) in the population that the unit represents. If the sampling weight is
54 ignored in the analysis, this is equivalent to setting all the weights to be equal to one,
55 producing biased estimates of population quantities such as means, totals and
56 proportions. For more discussion on survey weights see Levy & Lemeshow⁽²⁾, ch 16;
57 Valliant *et al.*⁽¹⁾ and Valliant & Dever⁽³⁾. Using the sampling weights but ignoring
58 the sampling design will result in biased estimates of the SEs associated with the
59 estimated population quantities, resulting in invalid inferences⁽⁴⁻⁶⁾.

60 A statistical agency may release unit-level data for public use with different levels of
61 confidentiality protection. There are essentially two ways the unit-level data, often
62 called the Confidentialised Unit Record File (CURF), are released: with or without
63 the sampling design variables. The purpose of the latter approach is to protect the
64 identity of the respondents. Instead of the sampling design variables, a set of replicate
65 weights are supplied; the number of which may vary from survey to survey. To obtain
66 unbiased estimates and valid estimates of SEs, the researcher needs to use either: (1)
67 the individual's sampling weight as well as the sampling design variables; or (2) the
68 individual's sampling weight in conjunction with the set of replicate weight variables
69 (Section 2.1). In the estimation of valid SEs, a Taylor series linearisation method is
70 applied in the first approach; and for the second approach, a replication method such
71 as the jackknife method is required. An example of the two procedures using the
72 NHANES data can be found in StataCorp⁽⁷⁾, p116–117. The importance of using the
73 sampling weights and the sample design variables as in (1) is demonstrated in Saylor
74 *et al.*⁽⁸⁾ and Kim *et al.*⁽⁹⁾ with particular reference to NHANES and the Korean
75 NHANES respectively.

76 Researchers new to survey analysis often struggle to understand the weighting
77 procedure and how this should be incorporated into the analysis. The focus of this
78 paper is to answer the following questions when the data supplied includes the
79 replicate weights rather than the sampling design variables, as in (2):

- 80 1. *What happens if I don't use the sampling weights or the design information in*
81 *my analysis?*;
- 82 2. *How do I carry out analyses such as estimation of means, proportions and their*
83 *SEs; and estimates of coefficients for a logistic regression model.*
- 84 3. *How do I obtain estimates for subgroups when data is sampled using a complex*
85 *survey design?*
- 86 4. *How do I set up the code to include to incorporate the replicate weights in Stata?*

87 The data from the Australian Health Survey (2011-13) (AHS) will be used to answer
88 these questions showing results for three analyses: (1) unweighted; (2) weighted but

89 not accounting for design; and (3) weighted and accounting for complex sample design.
90 This paper is structured as follows: in Section 2, a description of replicate weights, the
91 AHS sampling design and variables to be used are described, along with the details of
92 the statistical analyses performed. In Section 3, results for the three methods are
93 provided with discussion of results in Section 4.

94 **2 Experimental Methods**

95 **2.1 Replicate Weights**

96 In the replication approach, sub-samples are selected from the original sample in a way
97 that reflects the sample design and estimation methods used, analysis is carried out on
98 each sub-sample and the variance between these estimates is used to estimate the
99 variance and SE of the required parameter estimate from the full sample⁽¹⁰⁾. The
100 replicate weights may be constructed by the statistical agency using a jackknife
101 procedure which depends upon the sample design. Often in a multistage design, each
102 replicate includes all but one primary sampling unit (PSU) and the total number of
103 replicates is the number of PSUs in the design⁽⁶⁾. If the sample design involves a large
104 number of PSUs, there will be a large number of replicates. An alternative is the
105 delete-a -group jackknife method^(10,11) where each replicate is formed by deleting one
106 in R of the PSUs, where R is the number of replicates required. For more detail in
107 how to generate replicate weights (in Stata) given the sample design refer to Valliant
108 & Dever⁽³⁾, Section 5.4.

109 When the survey data set does not provide the survey design variables, the number of
110 PSUs (the top level cluster variable) is often not provided. Instead, the number of
111 replicate weights and the variable names will be specified in the user documentation.
112 When the statistical agency constructs and supplies the set of replicate weights in a
113 CURF, it simplifies the task for the analyst in as much as the variables pertaining to
114 the sampling design as in approach (1) and the syntax required in statistical software
115 to use them are not required. However, for approach (2), the data analyst must know
116 how to use the replicate weights, a demonstration of which is given in this paper.

117 For the AHS, the set of replicate weight consists of $R = 60$ variables, in addition to the
118 individual's sampling weight (referred to as the person weight). Each of these 60
119 replicate weight variables will have a collection of rows or individuals where the weight
120 is set to zero, such that no two variables will have the same rows set to zero but across
121 the 60 variables, each case will appear as zero in one variable only. The collection of
122 rows which are set to zero for a replicate weight variable indicate those individuals
123 that are deleted to form the replicate. Because replicates are formed by deleting PSUs,
124 the number of rows set to zero in each variable may vary. In each replicate weight
125 variable, the remaining non-zero weights are adjusted to sum to the number of units in
126 the population; so the sum of the weights for each of these variables is identical. It
127 would be incorrect to only use a subset of the full set of replicate weight variables.
128 This set of replicate weights is then used in the jackknife variance estimation for the
129 parameters of interest. For more details about standard errors and the replicate
130 weights technique for the AHS see see the AHS User's Guide⁽¹⁰⁾ and for an
131 introduction to jackknife estimation see Abdi & Williams⁽¹²⁾.

132 **2.2 Data Description**

133 The Australian Health Survey (2011-13) (AHS) combines three national health surveys
134 collected by the Australian Bureau of Statistics (ABS), namely: National Health
135 Survey (NHS); National Nutrition and Physical Activity Survey (NNPAS); and
136 National Health Measures Survey (NHMS) which is a biomedical information
137 component. Information collected includes health status, risk factors, actions and
138 socioeconomic circumstances. More detailed information about the structure of the
139 AHS may be found in the AHS First Results Report⁽¹³⁾.

140 For the purpose of this paper, variables used will be measures taken from NNPAS, as
141 this is the survey generally of most interest for nutrition related questions. The
142 sampling design used a stratified multistage area sample of private dwellings, collecting
143 information by face-to-face interview. The strata are Statistical Divisions within each
144 state and territory; each stratum comprises a number of Census Collection Districts
145 (CDs) which were used as PSUs and consists of an average 250 dwellings. The CDs

146 were sampled within each stratum and then dwellings within a sample of a selected
147 block in each selected CD were selected. A total of 3047 PSUs were selected ; persons
148 were then randomly selected from each dwelling such that one adult and one child
149 aged 2-17 years were selected where possible. Oversampling, (i.e. higher sampling
150 rate) of older adults (65+ years) was also carried out. More detail of the sampling
151 design may be found in the AHS Users' Guide⁽¹⁰⁾. This complex survey design is
152 typical of many national surveys. The total responding sample ($n = 12153$) comprised
153 both adults and children aged 2+ years, and our analysis has been limited to adults
154 (aged 18+ years; $n = 9435$).

155 The survey included the collection of measured height (in cm) and weight (in kg) and
156 then Body Mass Index (BMI) was calculated as the weight in kilograms divided by the
157 square of the height in metres. BMI values are categorised according to the World
158 Health Organisation (WHO) and the National Health and Medical research Council
159 (NHMRC) guidelines. These categories are: *underweight* (< 18.5), *normal*
160 ($18.50-24.99$), *overweight* (25.00 to 29.99) and *obese* ($30+$)⁽¹⁰⁾. The relevant original
161 variable names in the NNPAS CURF are: weight (*PHDKGWBC*), height
162 (*PHDCMHBC*), measured BMI (*BMISC*) and BMI categories (*BMICATHY*).

163 There are three types of sampling weights supplied in the NNPAS dataset: household
164 weight; two person weights (for all responding persons and biomedical sample only).
165 For estimating mean BMI and proportions of persons categorised as overweight or
166 obese, the person weight (*NPAFINWT*) applied to all responding persons is
167 appropriate. The 60 replicate weights are named *WPM0101* - *WPM0160*.

168 **2.3 Statistical Analysis**

169 Estimating descriptive statistics and their SEs for a mixture of variable types was
170 conducted: continuous variables *Height (in cm)*, *Weight (in kg)* and *BMI*; categorical
171 variables *Overweight or Obese* and *Current smoker*. Coefficients for a logistic
172 regression model for the binary variable of *Overweight or Obese* were also estimated.

173 Three methods of statistical analyses were conducted:

174 (1) Unweighted: without sampling weights or replicate weights;

- 175 (2) Weighted: with sampling weights but without accounting for the complex design;
 176 equivalent to weighted analysis assuming simple random sampling (SRS); and
- 177 (3) Complex design: with sampling weights and accounting for the complex design
 178 using a jackknife procedure with the replicate weights.

179 For the three continuous variables *Height (in cm)*, *Weight (in kg)* and *BMI*, the
 180 estimated mean and SE were determined. A binary variable identifying adults (≥ 18
 181 years) was first created from the continuous age variable (*AGEC*); then a binary
 182 variable identifying overweight or obese adults was created. For the categorical variable
 183 for smoking *SMOKEQ1*, the percentage of current smokers is estimated for the adult
 184 population. A logistic regression model for the status of overweight or obese adults is
 185 applied using the covariates: *sex*, *age (in years)*, highest year of school completed
 186 (*SchEd*), total minutes undertaken physical activity in last week (*PhysActMin*),
 187 remoteness of area category (*ARIABC*) and current smoker (*SMOKEQ1*). Reference
 188 categories: for *sex* is male; for *SchEd* is year 12 or equivalent; for *ARIABC* is major
 189 city; for *Smoker* is yes. Stata 15 was used for all analyses, the commands for **mean**,
 190 **proportion**, and **logistic** were used with the appropriate **svy** command settings for
 191 the three methods given in the appendix. The 60 jackknife replicate weight variables
 192 are defined in Stata with the **jkrweight** option in the **svyset** command.

193 The formula used for the three methods are shown here for the estimate of the
 194 population mean and its variance.

- (1) Unweighted: the familiar sample mean of a single variable y , denoted by \bar{y} , and its estimated variance assuming a simple random sample without replacement of size n from a target population of size N , and sample variance s^2 , is calculated without the sampling weights or the replicate weights. If y_i is the i th observation ($i = 1 \dots n$) from the sample, then the sample mean and estimated variance is given by:

$$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \quad v(\bar{y}) = \left(1 - \frac{n}{N}\right) \frac{s^2}{n} \quad s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$$

195 (2) Weighted: If the sampling weight for an individual in the sample is denoted by
 196 w_i ($i = 1 \dots n$) where the sum of the weights is equal to the population size N ,
 197 $\sum_{i=1}^n w_i = N$, the estimator of the population mean is the mean of the weighted
 198 observations; and the variance is the equivalent to weighted analysis assuming
 199 simple random sampling (SRS) such that:

$$\hat{\theta} = \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i} \quad v(\hat{\theta}) = \left(1 - \frac{n}{N}\right) \left(\frac{n}{n-1}\right) \frac{1}{N^2} \sum_{i=1}^n w_i^2 (y_i - \hat{\theta})^2$$

200 (3) Complex design: the sample weights are used to calculate the weighted mean as
 201 given in (2). The replicate weight variable for each replicate group is used to
 202 obtain each of the R replicate estimates of the mean resulting in $\hat{\theta}_1 \dots \hat{\theta}_R$. The
 203 variance estimate of $\hat{\theta}$ is then given by $v^*(\hat{\theta})$:

$$v^*(\hat{\theta}) = m \sum_{r=1}^R (\hat{\theta}_r - \hat{\theta})^2 \quad (1)$$

204 where the jackknife multiplier (m) is given by $m = \frac{R-1}{R}$. For the AHS data, a
 205 delete-a-group jackknife method of replicate weighting is used producing $R = 60$
 206 replicate weights, so $m = 59/60$ is the required multiplier⁽¹⁰⁾.

207 For the coefficients in a logistic regression, the variance of the unweighted estimates
 208 was estimated using standard methods. For weighted analysis, the variance ignoring
 209 the sample design was estimated using a linearization approach⁽¹⁴⁾. The jackknife
 210 approach uses $v^*(\hat{\theta})$ defined in (1) where $\hat{\theta}_r$ is the estimate obtained using the weights
 211 for replicate r .

212 2.4 Estimates for subgroups

213 Often a researcher is interested in estimating a quantity such as a mean or proportion
 214 for a subgroup of the population, for example, the mean *BMI* by *Sex* may be of
 215 interest. In this section, we focus on how to carry out such analysis when the jackknife
 216 replication method is to be used. The subgroup analyses for the mean *BMI* and the
 217 proportion of *Overweight or Obese* and *Smoker* by *Sex* are conducted.

218 When analysing survey data for particular subgroups, it is strongly recommended that
219 the full data set be used rather than restricting the data to the particular cases
220 belonging to the subgroup. The Stata manual for survey data⁽⁷⁾ describes the use of
221 command options `subpop` and `over` when estimating parameters for subgroups of the
222 population rather than restricting the number of cases using conditional `if` or `in`
223 qualifiers. The `subpop` option can be used to break down estimates into two groups
224 using a binary variable with zero/nonzero values such as 0/1. The `over` option allows
225 a breakdown by a categorical variable with two or more categories. West *et al.*⁽¹⁵⁾
226 explain the conceptual differences between the ‘conditional’ approach and the design
227 (or ‘unconditional’) approach. The latter allows for the variance arising from the
228 randomness of subgroup representation in the sample.

229 The Stata manual⁽⁷⁾ p72 states that ‘every (replicate) weight variable produces a
230 replicate, even if it does not contain an observation within the subpopulation specified
231 in the `subpop` option’. However, if the `if` or `in` qualifiers are used instead, only the
232 PSUs that have at least one observation within the subgroup will produce a replicate.
233 However, West *et al.*⁽¹⁵⁾, p527 state that a conditional approach for complex design
234 estimation of variances will produce the correct standard errors assuming the replicate
235 weights have been correctly produced according to the complex sampling design.

236 Valliant *et al.*⁽¹⁾, p421 note that the jackknife correctly handles subgroup estimation
237 without the need to explicitly give people not in the subgroup a zero response variable.
238 They also note that for linearization approach to variance estimation, restricting the
239 data set to the subgroup is generally a mistake.

240 **3 Results**

241 The results for the descriptive statistics for the 5 chosen variables are listed in Table 1.
242 The unweighted (1) estimates of the means and percentages differ from the weighted
243 and complex design, as expected. For *Height* the unweighted mean is lower but for the
244 other listed variables, it is higher. The standard errors across the three methods are
245 different: the difference being larger between Methods (1) and (2) than between

246 Methods (2) and (3). What is interesting is the difference between the SEs for
247 methods (2) and (3). For *Height*, *Weight* and *BMI*, the SEs have decreased, however,
248 increases in SE are evident for *Overweight or Obese* (from 0.74 to 0.78) and for
249 *Current Smoker* (0.52 to 0.53). Method (1) gives a biased estimate of both the mean
250 (or percentage) and the SE; Method (2) gives an unbiased estimate of the mean (or
251 percentage) but a biased estimate of the associated SE; and Method (3) provides an
252 unbiased estimate of both the mean (or percentage) and its SE.

253 The results for the logistic regression model of whether or not an adult is *Overweight*
254 *or Obese* are provided in Table 2. The odds ratio (OR), standard error (SE), the
255 t-statistic, related p-value and the 95% confidence interval (CI) are shown for each of
256 the six covariates for each of the three methods. Similarly to Table 1, the results for
257 Method (1) unweighted analysis provides a different parameter estimate for the odds
258 ratio; but for methods (2) and (3), the estimated OR are equal for the corresponding
259 covariate.

260 The estimated SEs for the corresponding covariates differ across the three methods.
261 Comparing the SEs between Method (1) and Method(2), it is clear that for Method
262 (2), the SEs are all higher than for Method (1). All but one of the SEs are higher for
263 the complex design (3) results than for the weighted (2) results; with the SE for
264 *Smoker* being the exception. The most notable difference is for the covariate *Smoker*.
265 For *Smoker*, the unweighted method gives an $OR=1.20$ which is statistically
266 significantly higher than 1.0 (assuming a 5% level) with $t=2.89$, $P=0.004$ and 95% CI:
267 $(1.061, 1.365)$. However, for the complex design method, the result is not significant
268 $t=0.71$, $P=0.480$ and 95% CI: $(0.894, 1.267)$, thus underlining that invalid inferences
269 can be made if analysis does not take the complex design into account. Also
270 noteworthy are the results for the variable for Remoteness of area category:
271 (*ARIABC*). Method (1) reports, for the *Other* category, $SE=0.0805$, $t=2.03$ and
272 $P=0.043$ whereas the corresponding results for Method (3) gives $SE=0.1460$, $t=2.47$
273 and $P=0.016$. The results for Method (2) are similar to those for Method (3) with SE
274 slightly higher for Method (3).

275 To summarise the different SEs between methods for the same covariate, the ratio of

276 the SE for Method (2) to the SE for Method (1) found a minimum ratio of 1.12 (for
277 *Smoker-No*), a maximum of 1.78 (for *ARIABC-Other*) and a median of 1.36 across the
278 covariates. The ratio of the SE for Method (3) to SE for Method (2) found a minimum
279 ratio of 0.99 (for *Smoker-No*), a maximum of 1.19 (for *SchEd-Year 10*) and a median
280 of 1.07.

281 The subgroups analysis for the mean *BMI* and proportion of *Overweight or Obese* and
282 *Smoker* by *Sex* were carried out for the jackknife variance estimates. However, since
283 the same SE is produced for complex design estimates whether or not the ‘conditional’
284 or ‘unconditional’ approach is adopted, the results are not shown here. It is suggested
285 that the analyst becomes familiar with, and uses the ‘unconditional’ approach for good
286 practice, whether or not replicate weights are supplied in the data⁽⁵⁾. If the survey
287 design variables are supplied rather than the replicate weights, then the conditional
288 approach will produce incorrect SEs, hence the ‘unconditional’ approach is required.
289 Stata code for this analysis is shown in the appendix.

290 4 Discussion

291 When reading the literature on secondary analyses of national health surveys, it is
292 sometimes not clear whether or not the reported estimates are the weighted estimates
293 and whether or not the analysis accounted for the complex survey design. Bell *et al.*⁽⁴⁾
294 carried out a review of 1003 published papers reporting empirical research from 1995
295 to 2010 in three health surveys. They found that ‘60% of articles reported accounting
296 for design effects and 61% reported using sample weights’. For an Australian example,
297 Allman-Farinelli *et al.*⁽¹⁶⁾ examine BMI and the prevalence of overweight and obesity
298 by occupation using National Health Survey (NHS) 2004-05 data collected by the
299 ABS. It is reported that the person sampling weight is used in the analysis but there is
300 no mention of the use of the method to obtain the reported standard errors that
301 account for the complex sample design and how the restriction to adults aged 20-64
302 years was handled. The AHS data from 2011-12 was used in a study on cardiovascular
303 health by Peng *et al.*⁽¹⁷⁾. Poisson and logistic regression analyses were conducted on a

304 restricted subgroup of the core sample with analysis applying the biomedical sample
305 weights and jackknife method as recommended by ABS⁽¹⁸⁾.

306 Saylor *et al.*⁽⁸⁾ demonstrate the importance of using the sampling weights and the
307 survey's complex sample design in any statistical analysis with particular reference to
308 NHANES. The sampling design variables for NHANES, including the stratification
309 and cluster variables, are supplied in the data files in addition to the sampling weight.
310 The authors undertake analyses in SPSS including descriptive statistics, linear and
311 logistic regression using three methods: unweighted, weighted and complex samples.
312 They conclude that accurate estimates of means and frequencies are produced if using
313 weights without the complex sample design information but that 'weighting alone
314 leads to inappropriate population estimates of variability'⁽⁸⁾ p236. Similarly, Kim
315 *et al.*⁽⁹⁾ report that only 19.8% of the 247 research articles the Korean National
316 Health and Nutrition Examination Survey (KNHANES) cited in PubMed from 2007
317 -2012 correctly used survey design analysis. Using SAS and SUDAAN, Kim *et al.*⁽⁹⁾
318 compare the estimates of levels of lead, cadmium and mercury in the blood and the
319 associated SEs as well as odds ratios (and 95% confidence intervals) for hypertension
320 and osteoporosis for particular subgroups using both unweighted and a weighted
321 analysis accounting for the complex design features. The results highlight the
322 differences in the parameter estimates if weighting is not applied and the tendency for
323 SEs to be underestimated and the CIs to be invalid.

324 The weighted SRS SE estimator [Method (2)] treats the data as a simple random
325 sample of weighted values. This estimator at least partially accounts for the use of
326 weights but does not reflect the effect of stratification and clustering in the sample
327 design or the use of post-stratification in the estimation. Ignoring the effect of
328 stratification will mean that the estimator will tend to overestimate the true SE, while
329 ignoring the clustering and post-stratification will tend to underestimate the SE. The
330 net effect of these factors will depend on the particular design used and the variable
331 being considered and means that it is possible for the weighted SRS SE estimator
332 [Method (2)] to be larger or smaller than the SE estimates obtained using the replicate
333 weights [(Method (3))], which properly account for these effects. The clustering in the

334 NHS is not high, with an average of less than 7 dwellings selected per PSU and so we
335 would not expect a large increase in SE due to the clustering in the sample, but some
336 increase is evident in the complex variance analysis. A larger SE will result in a wider
337 confidence interval and in turn will reduce the power of an analysis⁽¹⁹⁾.

338 4.1 Conclusion

339 This paper discusses the results of three approaches to secondary analysis of complex
340 survey data which have replicate weight variables supplied in the data rather than the
341 survey design variables, such as the variables indicating the strata and cluster to which
342 people belong. These are important considerations for nutrition related analyses in
343 surveys employing replicate weights.

344 The first question: *What happens if I don't use the sampling weights or the design*
345 *information in my analysis?* is answered in two parts. If the sampling weights are not
346 used in the analysis, biased estimates are produced which will lead to incorrect
347 conclusions. This is demonstrated by the differences in the estimates produced with
348 the unweighted and weighted methods. In addition, if the complex design is not
349 included, which translated to not using the replicate weight variables, the standard
350 errors will also be incorrect. The use of these incorrect standard errors may result in
351 incorrect inferences and conclusions.

352 The second question posed was: *How do I carry out analyses such as estimation of*
353 *means, proportions and their SEs; and estimates of coefficients for a logistic regression*
354 *model.* This paper demonstrates the use of replicate weights for analysing complex
355 survey data by way of several examples using data from the Australian Health Survey
356 (AHS). The AHS data has 60 replicate weight variables supplied with the data file.
357 Other researchers of AHS data or other surveys with replicate weights may use this
358 analysis as an example.

359 For the third question: *How do I obtain estimates for subgroups when data is sampled*
360 *using a complex survey design?* We found no difference when the conditional or
361 'unconditional' approach was adopted using the replicate weights; thereby showing
362 that the approach using the replicate weights is robust and simplifies procedures for

363 the analyst. However, for good practice, we suggest that the analyst becomes familiar
364 with the unconditional approach to the analysis of subgroups when using other
365 approaches to variance estimation.

366 The Stata code for all analyses is provided, which answers the last question posed:
367 *How do I set up the code to include to incorporate the replicate weights in Stata?* This
368 code may be used as an example for researchers performing similar analysis. It is
369 recommended that the analyst refers to the user's guide for the particular survey and
370 the relevant Stata documentation to determine the type of replication method carried
371 out and the number of replicate weights to apply. Some further examples may be
372 found in chapter 5 of Valliant & Dever⁽³⁾.

373 References

- 374 1. Valliant R, Dever JA, Kreuter F (2013) *Practical Tools for Designing and*
375 *Weighting Survey Samples*. Statistics for Social and Behavioral Sciences, New
376 York: Springer.
- 377 2. Levy PS, Lemeshow S (2008) *Sampling of Populations: Methods and Applications*.
378 Wiley Series in Survey Methodology, Hoboken, New Jersey: John Wiley & Sons,
379 4th edition.
- 380 3. Valliant R, Dever JA (2018) *Survey Weights - A Step-By-Step Guide to*
381 *Calculation*. College Station, Texas: Stata Press.
- 382 4. Bell BA, Onwuegbuzie AJ, Ferron JM, *et al.* (2012) Use of design effects and
383 sample weights in complex health survey data: A review of published articles using
384 data from 3 commonly used adolescent health surveys. *American Journal of Public*
385 *Health* **102**, 1399–1405.
- 386 5. Heeringa SG, West BT, Berglund PA (2010) *Applied Survey Data Analysis*.
387 Statistics in the Social and Behavioral Sciences Series, Boca Raton, FL, U.S.A.:
388 Chapman and Hall / CRC, Taylor and Francis Group.
- 389 6. Campbell RT, Berbaum ML (2010) *Handbook of Survey Research*, chapter
390 Analysis of Data from Complex Surveys. Bingley UK: Emerald Group Publishing
391 Limited, 2nd edition, pp. 221–259.
- 392 7. StataCorp (2017) *Stata Survey Data Reference Manual, Release 15*. StataCorp
393 LLC, College Station, Texas.
- 394 8. Saylor J, Friedmann E, Lee HJ (2012) Navigating complex sample analysis using
395 national survey data. *Nursing Research* **61**, 231–237.
- 396 9. Kim Y, Park S, Kim NS, *et al.* (2013) Inappropriate survey design analysis of the
397 Korean National Health and Nutrition Examination Survey may produce biased
398 results. *Journal of Preventive Medicine & Public Health* **46**, 96–104.

- 399 10. ABS (2012) Australian Health Survey: Users' Guide, 2011-13, Cat. no.
400 4363.0.55.001. Technical report, Australian Bureau of Statistics (ABS), Canberra,
401 accessed 4/12/2017, Available at
402 [http://www.abs.gov.au/AUSSTATS/abs@.nsf/Lookup/4363.0.55.001Main+Features12011-](http://www.abs.gov.au/AUSSTATS/abs@.nsf/Lookup/4363.0.55.001Main+Features12011-13?OpenDocument)
403 13?OpenDocument.
- 404 11. Kott PS (2001) The delete-a-group jackknife. *Journal of Official Statistics* **17**,
405 521–526.
- 406 12. Abdi H, Williams LJ (2010) *Encyclopedia of Research Design*, chapter Jackknife.
407 Thousand Oaks, CA: Sage.
- 408 13. ABS (2012) Australian Health Survey: First Results, 2011-12, Cat. no.
409 4364.0.55.001. Technical report, Australian Bureau of Statistics (ABS), Canberra,
410 accessed 4/12/2017. Available at
411 [http://www.abs.gov.au/ausstats/abs@.nsf/Lookup/4364.0.55.001main+features12011-](http://www.abs.gov.au/ausstats/abs@.nsf/Lookup/4364.0.55.001main+features12011-12)
412 12.
- 413 14. Binder DA (1983) On the variance of asymptotically normal estimators from
414 complex surveys. *International Statistical Review* **51**, 279–292.
- 415 15. West BT, Berglund P, Heeringa SG (2008) A closer examination of subpopulation
416 analysis of complex-sample survey data. *The Stata Journal* **8**, 520–531.
- 417 16. Allman-Farinelli MA, Chey T, Merom D, *et al.* (2010) Occupational risk of
418 overweight and obesity: an analysis of the Australian Health Survey. *Journal Of*
419 *Occupational Medicine and Toxicology* **5**. Published Online: DOI:
420 10.1186/1745-6673-5-14.
- 421 17. Peng Y, Wang Z, Dong B, *et al.* (2017) Life's simple 7 and ischemic heart disease
422 in the general Australian population. *PLoS ONE* **12**, 1–9. Available at
423 <https://doi.org/10.1371/journal.pone.0187020>.
- 424 18. ABS (2005) Technical Manual: National Aboriginal and Torres Strait Islander
425 Health Survey, Expanded CURF, 2004-05, Cat. no. 4715.0.55.002. Technical

426 report, Australian Bureau of Statistics (ABS), Canberra, accessed 1/2/2018.
427 Available at
428 <http://www.abs.gov.au/AUSSTATS/abs@.nsf/Latestproducts/4715.0.55.002Main>
429 Features3002004-05.

430 19. Burden S, Probst Y, Steel D, *et al.* (2012) The impact of complex survey design on
431 prevalence estimates of intakes of food groups in the Australian National
432 Children's Nutrition and Physical Activity Survey. *Public Health Nutrition* **15**,
433 1362–1372.

434 **A Stata code: Estimates of Means and Proportions**

435 The code in this section relates to the results in Table 1. In AHS data, some variables
436 have been given missing codes of 98, 99, 997, 998 and 999. These values were replaced
437 with appropriate codes for missing observations in Stata such as .a, .b, .c (Stata code
438 not given). For convenience, the variable for weight (*PHDKGWBC*) was renamed to
439 *Weight_kg*. Similarly, the variable for height (*PHDCMHBC*) was renamed to
440 *Height_cm*.

441 A dummy variable to indicate adults was created:

```
442 gen Adults=0 if AGE<18
443   replace Adults=1 if AGE<=18 & AGE<.
444   label define Adultslabel 0 "Under 18" 1 "18 or over"
445   label values Adults Adultslabel
```

446 A dummy variable to indicate the BMI category of overweight or obese was also
447 created:

```
448 gen OverObese = 0
449   replace OverObese = 1 if BMISC >= 25 & BMISC<.
450   replace OverObese = .a if BMISC==.a
451   replace OverObese = .b if BMISC==.b
452   label define OverOblabel 1 "Overweight or Obese" 0 "Not Overweight or
453   Obese" .a "Measurement not taken - refusal" .b "Measurement not taken -
454   other reason"
455   label values OverObese OverOblabel
456   codebook OverObese, tabulate(20)
```

457 **A.1 Method (1): Unweighted**

458 Unweighted results are obtained using standard procedures without sampling weights
459 or accounting for design features.

```
460 mean Height_cm, over(Adults)
```

```
461 mean Weight_kg, over(Adults)
462 mean BMISC, over(Adults)
463 proportion OverObese, over(Adults)
464 proportion SMOKEQ1, over(Adults)
```

465 Alternatively, the same results may be obtained using the following code assuming
466 simple random sampling:

```
467 svyset, clear
468 svyset _n

469 svy, subpop(Adults): mean Height_cm
470 svy, subpop(Adults): mean Weight_kg
471 svy, subpop(Adults): mean BMISC

472 svy, subpop(Adults): proportion OverObese
473 svy, subpop(Adults): proportion SMOKEQ1
```

474 A.2 Method (2): Weighted

475 Weighted results include the sampling weights but does not account for the complex
476 sample design. *NPAFINWT* are the sampling weights supplied with the data.

```
477 svyset, clear
478 svyset _n [pweight=NPAFINWT]

479 svy, subpop(Adults): mean Height_cm
480 svy, subpop(Adults): mean Weight_kg
481 svy, subpop(Adults): mean BMISC

482 svy, subpop(Adults): proportion OverObese
483 svy, subpop(Adults): proportion SMOKEQ1
```

484 **A.3 Method (3): Complex Design**

485 Results for Method (3) are weighted and account for the complex design: utilises the
486 sampling weights, *NPAFINWT*, and the 60 replicate weights *WPM0101 - WPM0160*
487 supplied with the data.

```
488 svyset, clear
489 local mult =59/60
490 svyset [pweight=NPAFINWT], jkrweight(WPM01*, multiplier('mult'))
491 vce(jackknife)

492 svy, subpop(Adults): mean Height_cm
493 svy, subpop(Adults): mean Weight_kg
494 svy, subpop(Adults): mean BMISC

495 svy, subpop(Adults): proportion OverObese
496 svy, subpop(Adults): proportion SMOKEQ1
```

497 **B Stata code: Logistic Regression**

498 The code in this section relates to the results in Table 2. For convenience, the
499 following variables were renamed: *LVHNSQBC* was renamed to *NonSchEd*;
500 *HYSCHCBC* was renamed to *SchEd*; and *EXLEVELN* was renamed to *PhysActMin*.

501 **B.1 Method (1): Unweighted**

```
502 svyset, clear
503 svyset _n

504 svy, subpop(Adults): logistic OverObese SEX AGECE PhysActMin i.SchEd
505 i.ARIABC i.SMOKEQ1
```

506 **B.2 Method (2): Weighted**

```
507 svyset, clear
```

```
508 svyset _n [pweight=NPAFINWT]
509 svy, subpop(Adults): logistic OverObese SEX AGECE PhysActMin i.SchEd
510 i.ARIABC i.SMOKEQ1
```

511 **B.3 Method (3): Complex Design**

```
512 svyset, clear
513 local mult =59/60
514 svyset [pweight=NPAFINWT], jkrweight(WPM01*, multiplier('mult'))
515 vce(jackknife)
516 svy, subpop(Adults): logistic OverObese SEX AGECE PhysActMin i.SchEd
517 i.ARIABC i.SMOKEQ1
```

Table 1: Results for estimates of mean Height (in cm), Weight (in kg), BMI, percentage of Overweight or Obese adults ($BMI \geq 25$) and percentage of Current Smokers for all adults with associated standard errors (SE) are shown for three methods: (1) Unweighted; (2) Weighted; and (3)Complex design.

	n	Parameter	(1)	(2)	(3)
		Estimate	Unweighted	Weighted	Complex design
Height (in cm)	8057	Mean	168.5297	169.1566	169.1566
		SE	0.1117	0.1518	0.1187
Weight (in kg)	8009	Mean	78.4144	78.0985	78.0985
		SE	0.1993	0.2686	0.2503
BMI	7958	Mean	27.5417	27.2245	27.2245
		SE	0.0624	0.0816	0.0801
Overweight/Obese	7958	Percentage	64.0990	61.9995	61.9995
		SE	0.5378	0.7438	0.7752
Smoker	9435	Percentage	18.9189	17.6150	17.6150
		SE	0.4032	0.5244	0.5287

Table 2: Results for logistic regression for all adults ($n=7874$): whether or not an adult is Overweight or Obese given 6 explanatory variables are shown for three methods: (1) Unweighted; (2) Weighted; and (3) Complex design. Odds ratio (OR) standard error (SE) the t-statistic, related p-value and 95% confidence interval (CI) are shown.

		OR	SE	t	P	95% CI	
(1) Unweighted							
Sex	F	0.5188	0.0257	-13.26	0.000	0.4708	0.5716
Age		1.0201	1.743E-3	11.67	0.000	1.0167	1.0236
PhysActMin		0.9996	0.833E-4	-5.04	0.000	0.9994	0.9997
SchEd	Year 11	1.4583	0.1247	4.41	0.000	1.2332	1.7244
	Year 10	1.3421	0.0877	4.50	0.000	1.1808	1.5254
	Year 9	1.4503	0.1624	3.32	0.001	1.1643	1.8064
	Year 8 or below	1.1046	0.1215	0.90	0.366	0.8904	1.3703
ARIABC	Inner regional	1.1491	0.0739	2.16	0.031	1.0130	1.3034
	Other	1.1523	0.0805	2.03	0.043	1.0047	1.3215
Smoker	No	1.2037	0.0773	2.89	0.004	1.0614	1.3650
Constant		1.5197	0.1775	3.58	0.000	1.2087	1.9108
(2) Weighted							
Sex	F	0.5091	0.0337	-10.19	0.000	0.4470	0.5797
Age		1.0262	2.378E-3	11.15	0.000	1.0215	1.0308
PhysActMin		0.9996	1.145E-4	-3.90	0.000	0.9993	0.9998
SchEd	Year 11	1.4171	0.1624	3.04	0.002	1.1320	1.7741
	Year 10	1.3097	0.1154	3.06	0.002	1.1020	1.5565
	Year 9	1.4918	0.2242	2.66	0.008	1.1111	2.0030
	Year 8 or below	1.2022	0.1760	1.26	0.208	0.9023	1.6017
ARIABC	Inner regional	1.3775	0.1203	3.67	0.000	1.1608	1.6346
	Other	1.3156	0.1429	2.52	0.012	1.0632	1.6278
Smoker	No	1.0640	0.0939	0.70	0.482	0.8951	1.2649
Constant		1.2413	0.1980	1.35	0.175	0.9080	1.6969
(3) Complex design							
Sex	F	0.5091	0.0371	-9.26	0.000	0.4399	0.5890
Age		1.0262	2.705E-3	9.80	0.000	1.0208	1.0316
PhysActMin		0.9996	1.174E-4	-3.80	0.000	0.99938	0.9998
SchEd	Year 11	1.4171	0.1834	2.69	0.009	1.0939	1.8360
	Year 10	1.3097	0.1372	2.58	0.013	1.0621	1.6150
	Year 9	1.4918	0.2377	2.51	0.015	1.0846	2.0520
	Year 8 or below	1.2022	0.1811	1.22	0.227	0.8893	1.6251
ARIABC	Inner regional	1.3775	0.1287	3.43	0.001	1.1426	1.6606
	Other	1.3156	0.1460	2.47	0.016	1.0536	1.6427
Smoker	No	1.0640	0.0929	0.71	0.480	0.8935	1.2671
Constant		1.2413	0.2217	1.21	0.231	0.8683	1.7745

Notes: Reference categories: for *Sex* is Male; for *SchEd* is Year 12 or equivalent; for *ARIABC* is Major city; for *Smoker* is Yes.