



Background

Like many large-scale health surveys, the Australian Longitudinal Study on Male Health (Ten to Men) used a complex sampling scheme. This choice was made because sampling the target population using a simple random sample was not feasible. Sampling theory therefore plays an important role in our study design because it provides a framework for efficiency gains [1]. In Ten to Men, the key elements of the sample design were the use of stratification, multi-stage sampling and cluster sampling to select prospective participants and invite them to take part in the study. This design has implications for the analysis of data from Ten to Men for both inferences about population means or prevalences, and for quantifying the magnitude of associations between exposures and outcomes. Such analysis implications are, however, often poorly understood. At the extreme, views differ on whether to *adjust* for aspects of the study design and sampling scheme at the analysis stage (including accounting for unequal sampling fractions using inverse-probability-of-selection sampling weights) or to *not adjust*. Korn and Graubard [2] give an excellent example of this controversy using US National Health and Nutrition Examination Surveys (NHANES). At the heart of this debate is a trade-off between miti-

the SA2 as the PSU), adjustment for stratification (no adjustment, adjustment using the stratification variable as a covariate, adjustment using the survey command), and use of sample weights (yes or no). We also examine the association between self-rated health and smoking status using logistic regression, where the effect size of interest is an odds ratio. We again omit the results from analyses that use a multi-level logistic model for the same reasons discussed in the previous section.

In an analysis that makes no adjustment for the multi-stage design or for stratification or weighting (Table 2, row A), the mean difference between the two groups is -5.1 kg (95 % CI -5.8 to -4.5 kg). That is, those who describe themselves as having very good or excellent health report are, on average, 5.1 kg lighter than those who have good, fair or poor health. Adjusting for stratification by using a series of indicator variables for remoteness to enter it into the model as a categorical variable (row B) also gives a mean difference of -5.1 kg with 95 % CI -5.7 to -4.4 kg. Repeating the analysis in row A but with the use of sample weights to adjust for bias gives a smaller difference of -4.4 kg, but with a wider confidence interval than observed previously (95 % CI -5.6 to -3.3). Adjustment for stratification makes only a small difference to this result (row D).

Repeating the analysis to account for all stages of sampling using a multilevel model (rows E and F) gives a mean difference of -4.9 kg (95 % CI -5.5 to -4.2), with further adjustment for stratification giving a difference of -4.8 kg (95 % CI -5.5 to -4.2). As with estimating population prevalences using multi-level models, it is not possible to easily account for the sample weighting in this context.

The final four rows in Table 2 show results obtained using the survey commands to estimate the population mean difference. When SA1s are defined as the PSU and sample weights are used (row G), the mean difference between the two groups is -4.4 kg (95 % CI -5.5 to -3.2). When no weights are used, the difference is -5.1 kg (95 % CI -5.8 to $-$

