
Subject: Pooled Datasets - Use of Svyset & regional controls

Posted by [lukassg](#) on Fri, 03 Jan 2014 13:48:01 GMT

[View Forum Message](#) <> [Reply to Message](#)

For a research paper we have access to 'Demographic and household surveys' from several different countries taken at different years, e.g. surveys for Uganda for the years 1996, 2002, 2006 and 2011. We then want to pool all surveys from one country to a single dataset.

We experience the following two challenges:

1)
In order to account for the complex survey design we think we have to correctly specify the weights, stratification and clusters for each survey. Even though each survey is from the same country, they can differ slightly depending on the year.
Thus when we pool them, we still want to correctly specify the survey design. However now the question arises how to do it. Before when doing each year by its own, we used code along the following line:

```
gen weight = v005 / 1000000  
egen stratid = group (v024 v025), label  
svyset [pweight=weight], psu(v021) strata(stratid)
```

The main thing that differs between the surveys is the stratification variables. Sometimes there exists already a stratification variable, sometimes we had to create one like above. Also sometimes the variable v024 (region) for example has 6 values in one year and 10 in the next year. Is it even possible to correctly stratify our dataset when we pool different surveys?

2)

Since we also want to control for regional / community effects later on in our regression models (using svy: reg or svy: logit/clogit) it can be problematic if the defined regions and clusters differ between the surveys.

The only solution we see, is performing single regressions for each year/survey. The drawback is that one cannot directly see whether differences in the constant term or the coefficient of maternal education between the different years/surveys are significant.

Is there any other statistical method that could deal with this dilemma?

Subject: Re: Pooled Datasets - Use of Svyset & regional controls

Posted by [Sarah B](#) on Mon, 03 Feb 2014 01:08:17 GMT

[View Forum Message](#) <> [Reply to Message](#)

Hi Lukas,

Most of these issues are addressed in a similar thread:
<http://userforum.measuredhs.com/index.php?t=tree&goto=93>

3&S=ab4524c1bf5aed945d825720798d7ea8#page_top

For Uganda specifically, there are a few extra concerns: regions have changed over time, and not all regions of the country were included in every survey.

If you want the same stratification in every survey, I believe you could regroup the strata into the 4 regions (central, eastern, northern, and western) shown on page 4 of the UDHS 2000-2001 final report. In each final report, I believe information is included that would allow you to re-categorize the data into these 4 groups, though you would have to do this by hand.

A downside of this approach is that the geographic regions are quite large, and thus would not represent well any kind of "community." You would also lose some of the efficiency gains from stratification in estimating your sampling errors. On the plus side, though, your regions would be comparable across surveys, and you could pool the data files together -- with the caveats about weights and renaming clusters and regions mentioned in the other thread.

By the way, the same clusters are not selected in different survey years -- a different sample is drawn for each survey. So the cluster identifiers should be renumbered/renamed to be unique within your pooled data file.

I hope that is helpful -- if you have further questions, please repost.
