Replication methods for analysis of complex survey data in Stata

North American Stata User's Group Meeting Boston August 2004

> Nicholas Winter Cornell University nw53@cornell.edu

Complex Survey Designs

- Stratification
- Clustering
- Unequal Probabilities of Selection
- → Traditional calculations give wrong point estimates; easily fixed through weighting
- → Traditional i.i.d. statistical calculations give wrong variance estimates

Approaches

Linearization

- Taylor Series Expansion for each statistic
- This is the world of Stata's svy
- Must be programmed separately for every estimator
- Requires information on stratum and PSU (ie, cluster) membership for each sample element

Replication Methods

- Take multiple "pseudo" samples (or replicates) from the dataset
- Variance calculated from the variance of the estimator across the replicates
- Once programmed, can plug in any estimator
- More and more public use datasets now include replicate weights

Replication Methods

- Balance Repeated Replication (BRR)
 - 2-PSU per stratum designs
 - Each replicate consists of half the PSUs (1 per stratum)
 - 2^L possible samples; can use appropriately selected subset
- Survey Jackknife
 - Drop one PSU from each replicate
 - Designs
 - 2-PSU per stratum (JK2)
 - 2+ PSU per stratum (JKn)
 - unstratified (JK1)

Basic Approach

- Given sampling weights and sample design information, calculate R sets of *replicate weights*
- Weights set to zero for excluded PSU(s)
- Other weights adjusted accordingly

Basic Approach (2)

 Variance of an estimator Θ is embarrassingly easy to calculate, once you have the replicate weights:

$$V(\Theta) = F \sum_{r=1}^{R} f_r \left(\Theta - \Theta_r\right)^2$$

$$V(\Theta) = F \sum_{r=1}^{R} f_r \left(\Theta - \Theta_r\right) \left(\Theta - \Theta_r\right)$$

across R replicates Θ is the full-sample estimate Θ_r is the estimate of Θ in the r'th replicate F is a technique-specific scaling factor f_r is a replicate-specific scaling factor (JKn only)

Advantages of Replication

- Easily extended to new techniques
 - No new programming for new estimators
- PSU and Stratum membership information may not be available
 - Privacy concerns are making replication more common on publiclyreleased datasets
- Easy to incorporate post-stratification or raking, and nonresponse adjustments into variance estimation
 - Simply apply the post-stratification (raking, NR adjustment) to each set of replicate weights in turn

Disadvantages

• Not implemented in Stata

Generating the weights

- survwgt (available on SSC)
 - Given sampling weight, strata, and PSU information
 - Calculates BRR, JK1, JK2, JKn replicate weights
 - Also does post-stratification, raking, nonresponse adjustments

Creating Replicate Weights

• Using National Health & Nutrition Examination Study (NHANES)

```
Stata Results
 use http://www.stata-press.com/data/r8/nhanes2d.dta, clear
 survwgt create brr , strata(strata) psu(psu) weight(finalwgt)
Obtaining hadamard matrix file...
Generating replicate weights.....
Created weights and set svr values:
   meth brr
     pw finalwgt
     rw brr_1 brr_2 brr_3 brr_4 brr_5 brr_6 brr_7 brr_8 brr_9 brr_10 brr_11
brr_12 brr_13 brr_14 brr_15 brr_16 brr_17 brr_18 brr_19 brr_20 brr_21
        brr_22 brr_23 brr_24 brr_25 brr_26 brr_27 brr_28 brr_29 brr_30 brr_31
        brr_32
    dof 31
    fay 0
   psun <not set>
```

Doing Analysis

- **svr** package (also available from SSC)
- Counterparts to official Stata's **svy** commands:
 - svrmean, svrtotal, svrratio
 - svrtab
 - svrmodel (for regression-style models: regression, logit/probit, ologit/oprobit, poisson, etc. etc.)

• And some extras

- svrcorr calculates variances for correlation coefficients
- svrest turns any command that accepts weights into a replication-based survey estimator (analogy to –simul- or –jknife-)

svrset

	svrset clear
	svrset set meth brr
•	svrset set pw finalwgt
•	svrset set rw brr_*
ł	svrset set dof 31
•	svrset set fay 0
•	<pre>svrset list meth brr pw finalwgt rw brr_1 brr_2 brr_3 brr_4 brr_5 brr_6 brr_7 brr_8 brr_9 brr_10 brr_11</pre>

svymean vs. svrmean

Stata Results										
. svymean bpsystol, by(sex)										
Survey mean est	imation									
pweight: finalı Strata: strata PSU: psu	vgt a		Num Num Pop	ber of obs ber of strata ber of PSUs ulation size	= 10351 = 31 = 62 = 1.172e+08					
Mean Subpop.	Estimate	Std. Err.	[95% Conf.	Interval]	Deff					
bpsystol Male Female	129.9253 124.2027	.6432933 .7051858	128.6132 122.7644	131.2373 125.6409	5.482304 5.162487					
. svrmean bpsyst Survey mean est	. svrmean bpsystol, by(sex) Survey mean estimation, replication (brr) variance method									
Analysis weight Replicate weight Number of replic k (Fay's method)	: finalwg ts: brr_1 cates: 32): 0.000	t.	Number Popula Degree	= 10351 = 1.172e+08 = 31						
Mean Subpop.	Estimate	Std. Err.	[95% Conf.	Interval]	Deff					
bpsystol Male Female	129.9253 124.2027	.6442178 .7094868	128.6114 122.7557	131.2391 125.6497	5.498073 5.225652					

svrtab [SVY] p. 78

svrtab race diabetes, row se ci fse(%5.4f) fci(%3.2f) fcell(%7.6f)cross-tabulation with replication-based (brr) standard errorsAnalysis weight:finalwgtAnalysis weight:finalwgtReplicate weights:brr_1Number of replicates:32Analysis method):0.000								
1=white, 2=black, 3=other	diab 0	etes, 1=yes, 0 1	=no Total					
White	0.968046 (0.0020) [0.96,0.97]	0.031954 (0.0020) [0.03,0.04]	1.000000					
Black	0.940965 (0.0062) [0.93,0.95]	0.059035 (0.0062) [0.05,0.07]	1.000000					
Other	0.979664 (0.0102) [0.94,0.99]	0.020336 (0.0102) [0.01,0.06]	1.000000					
Total	0.965754 (0.0018) [0.96,0.97]	0.034246 (0.0018) [0.03,0.04]	1.000000					
Key: row proportions (standard errors of row proportions) [95% confidence intervals for row proportions]								
Pearson: Uncorrected chi2(2) = 21.3483 Design-based F(1.52, 47.02) = 14.9413 P = 0.0000								

SVY] pp. 32-33

. svrmodel h	ighbp height v	weight age a	age2 fema	le black,	cmd(logit)	
Logit estimate	es with replic	cate-based ((brr) sta	ndard err	ors	
Analysis weig Replicate weig Number of rep k (Fay's metho	ht: fina ghts: brr_1 licates: 32 od): 0.000	lwgt L D		Number o Populati Degrees F(6, Prob > F	f obs = on size = of freedom = 26) = =	10351 1.172e+08 31 80.78 0.0000
highbp	Coef.	Std. Err.	t	P> t	[95% Conf.	Interval]
height weight age age2 female black _cons . lincom fema (1) female	0325996 .049074 .1541151 0010746 356497 .3429301 -4.89574 le+black + black = 0	.0058292 .003289 .0211487 .0002063 .0870778 .1484334 1.132855	-5.59 14.92 7.29 -5.21 -4.09 2.31 -4.32	0.000 0.000 0.000 0.000 0.028 0.000	0444884 .042366 .1109819 0014953 5340933 .0401982 -7.206214	0207109 .0557821 .1972482 0006539 1789007 .6456619 -2.585267
highbp	Coef.	Std. Err.	t	P> t	[95% Conf.	Interval]
(1)	0135669	.1727493	-0.08	0.938	3658915	.3387577

Postestimation

- svr routines work with the usual post-estimation commands
 - test (nee svytest)
 - lincom
 - etc.
- Some ugly programming here, but it works . . .

My personal favorite: svrest

. svrest "sum	height" "r(me	an) r(sd)"							
Estimates with	replication	(brr) base	d standar	d errors					
Command: sum height Analysis weight: finalwgt Replicate weights: brr_1 Number of replicates: 32 Degrees of freedom = 31									
	Coef.	Std. Err.	t	P> t	[95% Conf.	Interval]			
stats									
stat1 stat2	168.4599 9.699111	.1469785 .0754005	1146.15 128.63	0.000 0.000	168.1601 9.545331	168.7597 9.852891			

svrest (2)

. svi Estir	svrest "regress weight height sex" "e(r2)", matrices(e(b)) liststats									
Comma Analy	and: ysis weig	pht:	regres finalw	s weight l gt	height se	x				
Numbe	er of rep	plicates:	32			Degrees	of freedom =	31		
		C	oef.	Std. Err.	t	P> t	[95% Conf.	Interval]		
mat1	height sex _cons	.660 -3.90 -33.3	0245 1195 5467	.028589 .60011 5.58121	23.09 -6.50 -5.98	0.000 0.000 0.000	.6017168 -5.125127 -44.73762	.7183323 -2.677263 -21.97171		
stat	s stat1	. 261	7663	.0094207	27.79	0.000	. 2425527	.2809799		
Key:	height sex _cons stat1	height sex _cons e(r2)								

A Note on Development

- svrmodel is fairly simple, really
- Means, totals, ratios, and tabulation
 - Stata's **svy** commands are implemented as ado files
 - internal _svy calculates variances of means, totals & ratios
 - called by svymean, svytotal, svyratio and svytab
 - svrcalc.ado is reverse-engineered to return results in the form that _svy does, as expected by svymean, svytotal, svytab, etc.