Standard Error of Estimates in Complex Surveys: Estimating Village Household Milk Production at a sub-district level

Using Stata

Subir Mitra, B.E.(Electrical), PGDRM(IRMA)

Senior Manager (Sectoral Analysis & Studies), National Dairy Development Board, Anand

Aug 1, 2013



Outline

- One Stage Stratified Cluster Sampling Design, utilising GIS
- Introduction to the dataset
- Getting results using EXCEL & STATA
- Why Stata is better & quicker
- Acknowledging Statalist Forum
- Looking ahead!

Using GIS tool for areal stratification

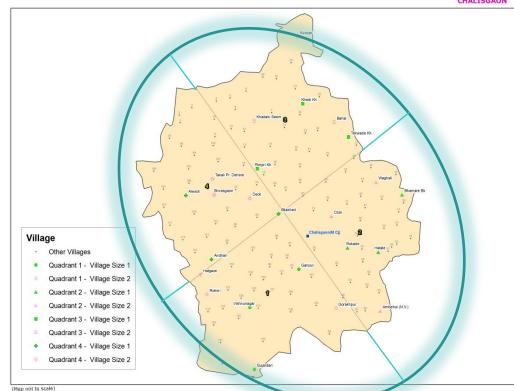


First, we use the **Directional Ellipse tool** of GIS to give geographic or areal representation, through GIS (ESRI ArcGIS) using digital maps of the sub-district plotted with village centroids.

Next, intersection of minor & major axis of the ellipse, gives us **4 geographical quadrants** within the ellipse, which we use as the first level of stratification.

Lastly, we stratify each of these 4 geographical strata into 2 further strata by grouping those villages which are higher & lower than the average number of households among all the villages in that quadrant.

Thus, now, we have **8 strata** - From each strata , we choose **2 sample villages** (clusters) randomly AND interview **all households** in those villages! In statistical jargon , this is called **Stratified One Stage Cluster Sampling**!



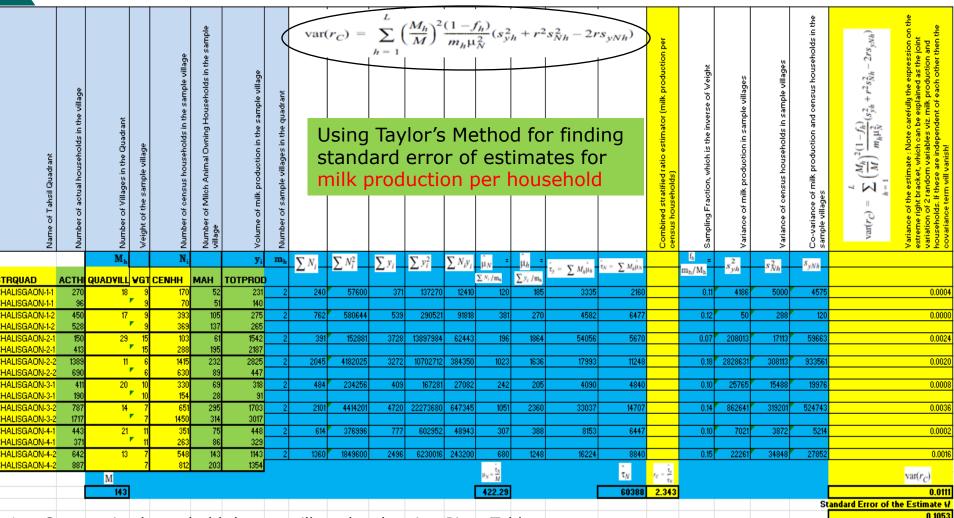
GEOQUAD	Average Households per Village	STRATA		Number of Households	SAMPLE VILLAGE	Number of Households
		CHALISGAON-1-1			Ganpur	170
		CHALISGAON-1-1	18	2762	Gujardari	70
		CHALISGAON-1-2			Hatgaon	369
1	288	CHALISGAON-1-2	17	7333	Rohini	393
		CHALISGAON-2-1			Bhamare Bk	288
		CHALISGAON-2-1	29	5160	Rokade	103
		CHALISGAON-2-2			Ozar	630
2	356	CHALISGAON-2-2	11	9062	Waghali	1415
		CHALISGAON-3-1			Khedi Kh	330
		CHALISGAON-3-1	20	4334	Tekwade Kh.	154
		CHALISGAON-3-2			Bahal	1450
3	441	CHALISGAON-3-2	14	10671	Kunzar	651
		CHALISGAON-4-1			Andhari	263
		CHALISGAON-4-1	21	4274	Bilakhed	351
		CHALISGAON-4-2			Deoli	812
4	376	CHALISGAON-4-2	13	8508	Shirasgaon	548
		TOTAL	143	52104		

Data Definitions & Some Village Household data samples.....

	SN	NAME OF VARIABLE	DEFINITION	REMARKS
ſ	1	TAHSIL	Name of tahsil	
	2	VILLAGE	Name of Village	
	3	SDTQVS	State District Tahsil Quadrant Village Sample Category	Quadrant Identification Number (each taluka cut into 4 quadrants)*
	4	VILL CD	Village Code (16 digit Census 2001)	
	5	HHNO	Household Number	
	6	FMLYMEM	Number of family members in the household	
	7	OPERLAND	Operation Land (acres)	
	8	CASTE	Caste	General=1, SC=2, ST=3, OBC=4
	9	OCCU	Occupation	Agri.=1, Agri. Labourer=2, Dairying=3, Service=4, Business=5, Others=6
4	10	ECOCATG	Economic Category	Above Poverty Line (APL)=1, Below Poverty Line (BPL)=2
	11	LCIM	Number of Local Cow In Milk	
	12	LCDRY	Number of Local Cow Dry	
	13	LCPROD	Milk Production (In Ltrs) of Local Cows	
	14	CBIM	Number of Crossbred Cow In Milk	
	15	CBDRY	Number of Crossbred Cow Dry	
ı	16	CBPROD	Milk Production (In Ltrs) of Crossbred Cows	
ı	17	BFIM	Number of Buffalo In Milk	
ı	18	BFDRY	Number of Buffalo Dry	
-	19	BFPROD	Milk Production (In Ltrs) of Buffaloes	
ı	20	TOTPROD	Total Milk Production (In Ltrs)	
-	21	PUR	Purchase of Milk (In Ltrs)	
-	22	CONS	Consumption by the Household (Ltrs)	
ı	23	SLDCS	Sale to DCS	
ı	24	SLLOCAL	Sale to Local consumers within village	
-	25	SLDUDHIA	Sale to Dudhia	
		SLOUTVILL	Sale outside village	
- 1	27	SLPVTDAIRY	Sale to Private Dairv	

IAHSIL	VILLAGE	VILLEGE	STRQUAD	HHINO	FIVILTIVICIVI	OFERLAND	CASTE	0000	ECOCATO	CONVI	LCDKI	LCFROD	CBIIVI	CDDKI	CBFROD	DEIIVI	DEDINI	BEEROD	TOTEROD	FUN	CONS	3LDC3	SELUCAL	SEDUDINA	SECOTALE	SEFVIDAINT
CHALISGAON	Bahal	27030013002751	CHALISGAON-3-2	1	10	17	4	1	1	0	0	0	7	0	40	0	0	0	40	0	5	35	0	0	0	0
CHALISGAON	Bahal	27030013002751	CHALISGAON-3-2	2	7	6	4	1	1	2	1	4.5	0	0	0	0	0	0	4.5	0	4.5	0	0	0	0	0
CHALISGAON	Bahal	27030013002751	CHALISGAON-3-2	3	4	5.5	4	1	1	0	0	0	0	0	0	1	0	4	4	0	2	0	0	2	0	0
CHALISGAON	Bahal	27030013002751	CHALISGAON-3-2	4	3	6.25	4	1	1	0	0	0	0	0	0	3	0	20	20	0	2	0	0	18	0	0
CHALISGAON	Bahal	27030013002751	CHALISGAON-3-2	5	10	14	4	1	1	1	0	3	0	0	0	1	0	4	7	0	1	0	0	6	0	0
CHALISGAON	Bahal	27030013002751	CHALISGAON-3-2	6	4	3	4	1	1	0	0	0	0	0	0	0	0	0	0	0.5	0.5	0	0	0	0	0
CHALISGAON	Bahal	27030013002751	CHALISGAON-3-2	7	5	5	4	1	1	0	0	0	1	0	8	0	0	0	8	0	1	0	0	7	0	0
CHALISGAON	Bahal	27030013002751	CHALISGAON-3-2	8	5	0	4	1	1	0	0	0	0	0	0	0	0	0	0	1	1	0	0	0	0	0
CHALISGAON	Bahal	27030013002751	CHALISGAON-3-2	9	4	13	4	1	1	0	0	0	0	0	0	1	0	7	7	0	1	0	1	5	0	0
CHALISGAON	Bahal	27030013002751	CHALISGAON-3-2	10	4	3	4	1	1	0	0	0	0	0	0	0	0	0	0	1	1	0	0	0	0	0

Getting Results Using Excel



- Summarise household data at village level, using Pivot Table
- 2. Feed the summarised data in the worksheet as above, carefully!
- 3. Repeat the process for all sub-districts, within one district, for one parameter of interest, at one time!
- 4. What if you have many such districts and couple of parameters to be estimated for each?

Getting Quicker & Better Results Using Stata!

```
use hh.dta
gen LC = LCIM+LCDRY
gen CB=CBIM+CBDRY
gen BF=BFIM+BFDRY
gen MILCH=LC+CB+BF
gen INMILK=LCIM+CBIM+BFIM
gen SALE= SLDCS+SLLOCAL+SLDUDHIA+SLOUTVILL+SLPVTDAIRY
gen MAH=1 if MILCH>0
replace MAH=0 if MAH==.
gen NMAH=1 if MILCH==0
replace NMAH=0 if NMAH==.
gen MPH=1 if TOTPROD>0
gen MSH=1 if SALE>0
gen CONSMAH=CONS if MAH==1
gen CONSNMAH=CONS if MAH==0
gen SURPLUS = TOTPROD - CONS
replace CASTE=5 if CASTE <1
replace OCCU=7 if OCCU<1
replace ECOCATG=4 if ECOCATG<1
tab CASTE, gen(CAST)
tab OCCU, gen(OCC)
tab ECOCATG, gen(ECO)
rename CAST1 GEN
rename OCC1 AGRI
rename ECO1 API.
collapse (first) STRQUAD (first) VILLCD (count) ACTHH=HHNO (sum) GEN (sum),......, by (VILLAGE)
save villraw.dta,replace
merge 1:1 VILLCD using base.dta
drop _merge
sort STRQUAD
save vill.dta, replace
```

gen pstrata=1

svyset VILLAGE [pw=WGT], strata(STRQUAD) poststrata(pstrata) postweight(THSLHH) fpc(QUADVILL)

svy:ratio MILCH CENHH /* estimated milch animal holding in the tahsil per hh*/

svy:ratio TOTPROD CENHH /* estimated milk producing HHs as a ratio of total HHs in the tahsil */

svy:ratio SURMAH CENHH /* estimated milk production in the tahsil per hh*/

svy:ratio SURPLUS CENHH /* estimated producer's milk surplus in the tahsil per hh*/

svy:ratio CONS TOTPROD /* estimated net milk consumption all HHs to production ratio in tahsil*/

svy:ratio CONSMAH TOTPROD /* estimated milk consumption of producing HHs to production ratio in tahsil*/



Estimating population parameters for the sub-district from the village level data



Comparing results between Excel & Stata

EXCEL

We can observe that the milk production per census household-> Standard Error of estimate->	2.3427 Ltr/day 0.1053 Ltr/day
Therefore, assuming normal distribution, for a 95% confidence interval	
the upper limit is	2.5491 Ltr/day
and the lower limit is	2.1362 Ltr/day

STATA

```
. svy:ratio TOTPROD CENHH
(running ratio on estimation sample)
Survey: Ratio estimation
                                     Number of obs
Number of strata =
                                                              16
Number of PSUs
                        16
                                     Population size =
                                                           52104
N. of poststrata =
                                     Design df
     ratio 1: TOTPROD/CENHH
                            Linearized
                             Std. Err.
                    Ratio
                                           [95% Conf. Interval]
                 2.342667
                             .1053275
                                           2.099782
    ratio 1
                                                        2.585553
```

The estimate and the standard error is the same, only the width of the limits are higher in Stata since it assumes a t-distribution, which is more accurate statistically speaking, as we have a very small sample size!!

Discovering & Using Stata: Acknowledging excellent support from Statalist Forum

```
From:
                       own er-statalist@h soh sun2. harvard e du on behalf of Stas Kolen ikov ⊸skolen ko@oma Loom ⊳
Sent:
                       Wednesday, January 30, 2013 8:00 PM
Torr
                       statalist(0) hap haun2.harvard.edu
Subject:
                       Re: st: one stage cluster with prelitiminary stratification
*** assuming n1, n2, m1, m2, N are contained in the identically named scalars.
gen wgt = scalar(n1)/10 if stratum==1
replace wgt = scalar(n2)/8 if stratum==2 assert !missing(wgt)

    cluster size does not matter.

ese option 1: poststrata.
gen pstrata = 1
gen popsize = scalar(N)
syvset cluster [pw=wgt], strata( stratum ) poststrata( pstrata ) postweight( popsize )
*** option 2: rescale weights:
saura west.
generate wgt2 = wgt*scalar(N)/r(sum)
syvset cluster [pw=wgt2], strata(stratum)

    Star Kolenikov, PhD, PStat (SSC) :: http://star.kolenikov.name

    Senior Survey Statistician, Abt SRBI :: work email kolenikovs at arbi dot com

    Opinions stated in this email are mine only, and do not reflect the position of my employer

On Wed, Jan 30, 2013 at 6:32 AM, Subir Mitra < subir@nddb.coop> wrote:

    ONE STAGE CLUSTER WITH PRELIMINARY STRATIFICATION -I stratify population N (members)

living in clusters, which is known) into 2 strata and randomly pick up 10 clusters from 1st
stratum and 8 clusters from 2nd stratum (stratum population n1 & n2 and total clusters m1 & m2
in both stratum also known) and all members of the clusters are sampled.
> Any guidance to me to find the syyset command in this case, assuming N, m 1,m2,n1 and n2
```

known and I want to make use of it? (The problem is from Schaeffer et al. 1996-328 problem 8.19)



Summary

Stata is very user friendly and much easy to learn. Inspiration from World Bank household survey analyses!

- Excellent support exists from practitioners.
- We look foreword to work together with Stata Users and exchange ideas!
- We are exploring using Stata in Monitoring & Evaluation- Baseline Survey & Follow-ups (National Dairy Plan 2012-17)

Propensity Score Matching or Regression?