

36-303: Sampling, Surveys and Society

Variance Calculations for Weights

Brian W. Junker
132E Baker Hall
brian@stat.cmu.edu

05 April 2012

1

Handouts & Announcements

- These Lecture Notes
- R Handout
- HW06 is online
 - This really is the last hw!
 - Due next Thu Apr 12
 - Contains updated list of due dates for rest of semester
- Last Midterm Exam Apr 17
 - Review Apr 12

05 April 2012

2

Outline

- Variance Calculations for Weights
 - Taylor Series
 - Random Partition
 - Jackknife

05 April 2012

3

Variance Calculations for Weights

- Most survey sample estimates have a ratio form:
$$\bar{y}_w = \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i}$$
- Two approaches to $Var(\bar{y}_w)$:
 - Use a ***one-term Taylor approximation*** to “linearize” the survey estimate, and apply CLT.
 - Use a ***replication scheme*** to create “replicate samples” by resampling the real sample and look at the variability among the replicates.
 - Non-overlapping replicates: E.g., *Random Partitions*
 - Overlapping replicates: E.g., *Jackknife Method*

05 April 2012

4

Taylor Series Approximation (Bkgd)

- The **Delta Method**

- We know that if

$$\hat{\theta} - \theta \sim N(0, \sigma^2/n)$$

then

$$a(\hat{\theta} - \theta) \sim N(0, a^2\sigma^2/n)$$

- We can extend this to a nonlinear function

$$f(\hat{\theta}) - f(\theta) = f'(\theta)(\hat{\theta} - \theta) + (\text{remainder})$$

so that

$$f(\hat{\theta}) - f(\theta) \approx f'(\theta)(\hat{\theta} - \theta) \sim N(0, [f'(\theta)]^2\sigma^2/n)$$

Taylor Series Approximation (Bkgd)

- Univariate Delta Method

$$\begin{aligned} \text{If } & \hat{\theta} - \theta \sim N(0, \sigma^2/n) \\ \text{then } & f(\hat{\theta}) - f(\theta) \sim N(0, [f'(\theta)]^2\sigma^2/n) \end{aligned}$$

- Multivariate Delta Method

$$\text{If } \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix} - \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \frac{1}{n} \Sigma \right)$$

then

$$\begin{aligned} f \begin{pmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \end{pmatrix} - f \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} \\ \sim N \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \frac{1}{n} \left(\frac{\partial f}{\partial \theta_1}, \frac{\partial f}{\partial \theta_2} \right) \Sigma \begin{pmatrix} \frac{\partial f}{\partial \theta_1} \\ \frac{\partial f}{\partial \theta_2} \end{pmatrix} \right) \end{aligned}$$

Taylor Series for Ratio Estimator

- Now we consider

$$\bar{y}_w = \frac{\sum_{i=1}^n w_i y_i}{\sum_{i=1}^n w_i} = \frac{\hat{\theta}_1}{\hat{\theta}_2} = f(\hat{\theta}_1, \hat{\theta}_2)$$

- The gradient of f has components

$$\frac{\partial f}{\partial \theta_1} = 1/\theta_2, \quad \frac{\partial f}{\partial \theta_2} = -\theta_1/\theta_2^2$$

- The Variance/Covariance Matrix for (θ_1, θ_2) is

$$\Sigma = \begin{bmatrix} \text{Var}(\sum_i w_i y_i) & \text{Cov}(\sum_i w_i y_i, \sum_i w_i) \\ \text{Cov}(\sum_i w_i y_i, \sum_i w_i) & \text{Var}(\sum_i w_i) \end{bmatrix}$$

Taylor Series Variance for Ratio Estimator

- Applying the Multivariate Delta Method we get

$$\begin{aligned} \text{Var}_{TS}(\bar{y}_w) \approx \\ \frac{1}{(\sum_i w_i)^2} \left[\text{Var}(\sum_i w_i y_i) - 2\bar{y}_w \text{Cov}(\sum_i w_i y_i, \sum_i w_i) + (\bar{y}_w)^2 \text{Var}(\sum_i w_i) \right] \end{aligned}$$

- Need to calculate the variances and covariance above – see next slide...

Calculating the Variances for TS Method...

If we assume that each pair $(w_i y_i, w_i)$ is independent of every other pair (not quite true but close!) then

$$\text{Var}\left(\sum_{i=1}^n w_i\right) = \sum_{i=1}^n \text{Var}(w_i) = n \text{Var}(w) \approx n \cdot \frac{1}{n-1} \sum_{i=1}^n (w_i - \bar{w})^2 = n \cdot s_w^2$$

where $\bar{w} = \frac{1}{n} \sum_i w_i$. Similarly,

$$\text{Var}\left(\sum_{i=1}^n w_i y_i\right) \approx n \cdot \frac{1}{n-1} \sum_{i=1}^n (w_i y_i - \overline{w y})^2 = n \cdot s_{wy}^2$$

where $\overline{w y} = \frac{1}{n} \sum_i w_i y_i$, and

$$\text{Cov}\left(\sum_{i=1}^n w_i y_i, \sum_{i=1}^n w_i\right) \approx n \cdot \frac{1}{n-1} \sum_{i=1}^n (w_i y_i - \overline{w y})(w_i - \bar{w}) = n \cdot s_{wy,w}$$

Example: HSS Advising Survey...

| Post-Strat. | Adv'ing OK | Samp Total | Prop | Pop Total | Prop | Weights |
|--------------|------------|------------|-------|-----------|-------|---------|
| Economics | 28 | 40 | 0.132 | 126 | 0.128 | 0.97 |
| English | 23 | 39 | 0.128 | 115 | 0.117 | 0.91 |
| History | 10 | 21 | 0.069 | 48 | 0.049 | 0.70 |
| ModLang | 3 | 8 | 0.026 | 16 | 0.016 | 0.62 |
| Philosophy | 1 | 4 | 0.013 | 7 | 0.007 | 0.54 |
| Psychology | 11 | 37 | 0.122 | 104 | 0.105 | 0.87 |
| SDS | 22 | 54 | 0.178 | 161 | 0.163 | 0.92 |
| Statistics | 3 | 6 | 0.020 | 8 | 0.008 | 0.41 |
| Interdisc/IS | 46 | 76 | 0.250 | 233 | 0.236 | 0.95 |
| Undeclared | 13 | 19 | 0.062 | 168 | 0.170 | 2.73 |
| Total | 160 | 304 | | 986 | | |

weight = (Population Proportion) / (Sample Proportion)

TS Variance Estimate, HSS Advising

Data...

$$\begin{aligned} y_i &= 1 \text{ (yes) or } 0 \text{ (no)} \\ \bar{y}_w &= 0.5507865 \\ \bar{w} &= 1.001678 \\ \overline{w y} &= 0.5517105 \end{aligned}$$

$$\text{Var}\left(\sum_i w_i\right) = n \cdot s_w^2 = (304)(0.2124) = 64.57$$

$$\text{Var}\left(\sum_i w_i y_i\right) = n \cdot s_{wy}^2 = (304)(0.4127) = 125.47$$

$$\text{Cov}\left(\sum_i w_i y_i, \sum_i w_i\right) = n \cdot s_{wy,w} = (304)(0.1637) = 49.75$$

So

$$\text{Var}_{TS}(\bar{y}_w) = (125.47 - 2(0.5507)(49.75) + (0.5507)^2 * (64.57)) / (304 * 1.0017)^2 = 0.000973$$

This is larger (typical!) than the naive variance based on $\hat{p} = \bar{y}$:

$$\hat{p}(1 - \hat{p})/n = (0.53)(1 - 0.53)/(304) = 0.000819$$

We should also multiply by the fpc = $1 - 304/986 = 0.69!$

Replication Scheme: Random Partitions

- We partition the data into $r = 1, \dots, c$ sub-samples, and calculate the weighted mean from each sub-sample
- Requirements of the sub-samples:
 - They are non-overlapping (disjoint subsets of the sample);
 - Their union is all of the original sample;
 - Each sub-sample should take observations from every stratum
- From each sub-sample we recalculate

$$\bar{y}_w^{(r)} = \frac{\sum_{i=1}^n w_i^{(r)} y_i^{(r)}}{\sum_{i=1}^n w_i^{(r)}}$$

- Note that the weights have to be recalculated each time as well!

Replication Scheme: Random Partitions

- This leads to a new estimate of the mean

$$\bar{y}_{rep} = \frac{1}{c} \sum_{r=1}^c \bar{y}_w^{(r)}$$

- and a simple estimate of the variance

$$Var(\bar{y}_{rep}) = \frac{1}{c} \left[\frac{1}{c-1} \sum_{r=1}^c (\bar{y}_w^{(r)} - \bar{y}_{rep})^2 \right]$$

- Takes a lot of computation but is straightforward to do, with a little programming!
- No example – look at Jackknife instead!

Replication Scheme: Jackknife

- From the original sample we create $r=1, 2, \dots, n$ *Jackknife samples* (of size $n-1$), by deleting one observation at a time from the original data.

- From each jackknife sample

- Recalculate the weights
- Recalculate

$$\bar{y}_w^{(r)} = \frac{\sum_{i=1}^n w_i^{(r)} y_i^{(r)}}{\sum_{i=1}^n w_i^{(r)}}$$

- Now calculate

$$\bar{y}_{JK} = \frac{1}{n} \sum_{r=1}^n \bar{y}_w^{(r)} \quad Var_{JK}(\bar{y}_w) = \frac{n-1}{n} \sum_{r=1}^n (\bar{y}_w^{(r)} - \bar{y}_{JK})^2$$

Example: HSS Advising Data (Again)

| Post-Strat. | Adv'ing | | Samp | | Pop | | Weights |
|--------------|------------|------------|-------|------------|-------|------|---------|
| | OK | Total | Prop | Total | Prop | | |
| Economics | 28 | 40 | 0.132 | 126 | 0.128 | 0.97 | |
| English | 23 | 39 | 0.128 | 115 | 0.117 | 0.91 | |
| History | 10 | 21 | 0.069 | 48 | 0.049 | 0.70 | |
| ModLang | 3 | 8 | 0.026 | 16 | 0.016 | 0.62 | |
| Philosophy | 1 | 4 | 0.013 | 7 | 0.007 | 0.54 | |
| Psychology | 11 | 37 | 0.122 | 104 | 0.105 | 0.87 | |
| SDS | 22 | 54 | 0.178 | 161 | 0.163 | 0.92 | |
| Statistics | 3 | 6 | 0.020 | 8 | 0.008 | 0.41 | |
| Interdisc/IS | 46 | 76 | 0.250 | 233 | 0.236 | 0.95 | |
| Undeclared | 13 | 19 | 0.062 | 168 | 0.170 | 2.73 | |
| Total | 160 | 304 | | 986 | | | |

~~weight = (Population Proportion) / (Sample Proportion)~~

JK Variance Estimate, HSS Advising Data...

- There are 304 Jackknife samples, of size 303 each.

- 28 jackknife samples omit one of the Econ 'yes' obs's
- 12 jackknife samples omit one of the Econ 'no' obs's
- 23 jackknife samples omit one of the English 'yes' obs's
- 16 jackknife samples omit one of the English 'no' obs's
- etc., etc. for the other 8 post-strata

- Calculate $\bar{y}_w^{(r)}$'s

- The first few unique $\bar{y}_w^{(r)}$ are

0.5478, 0.5488, 0.5490, 0.5493, 0.5495 ...

- (there are many duplicates!)

JK Variance Estimate, Continued

- Now we calculate

$$\bar{y}_{JK} = \frac{1}{304} \sum_{r=1}^{304} \bar{y}_w^{(r)} = 0.5508 \quad (= \bar{y}_w)$$

and

$$Var_{JK}(\bar{y}_w) = \frac{304 - 1}{304} \sum_{r=1}^{304} (\bar{y}_w^{(r)} - \bar{y}_{JK})^2 = 0.000963$$

- Very similar to TS Variance estimate:

$$Var_{TS}(\bar{y}_w) = 0.000973$$

Actual Calculations...

- See R handout... (is there someone in every group that knows a little R?)
- My recommendation:
 - If you know the formula, **Taylor Series** approx is really easy to carry out. However, for a new statistic, have to re-apply Delta Method.
 - **Jackknife** is harder to set up, but once it's done, it works for **all** possible statistics, not just weighted averages
 - As sample size grows, TS and JK produce same answers
 - (again, we should multiply by fpc = $(1 - \text{samp}) / (\text{pop})$)

Making a Confidence Interval

- Approx 95% confidence interval, based on the Jackknife standard error:

$$(0.5508 - 2 * \sqrt{(1 - 304/986)(0.000963)} \quad , \quad 0.5508 + 2 * \sqrt{(1 - 304/986)(0.000963)}) \\ (0.4992 \quad , \quad 0.6024)$$

- In our fictional example we know the true population proportion:

$$p_{pop} = 546/986 = 0.553$$

- We capture the true mean in this case

Review

- Variance Calculations for Weights
 - Taylor Series
 - Random Partition
 - Jackknife
- HW06 is online
 - This really is the last hw!
 - Due next Thu Apr 12
 - Contains updated list of due dates for rest of semester
- Last Midterm Exam Apr 17
 - Review Apr 12