

Stat 9100.3: Analysis of Complex Survey Data

1 Logistics

Instructor: Stas Kolenikov, kolenikovs@missouri.edu

Class period: MWF 1-1:50pm

Office hours: Middlebush 307A, Mon 1-2pm, Tue 1-2 pm, Thu 9-10am.

Website: Blackboard <http://courses.missouri.edu>

Information: This course covers some topics in modern analytical tools developed for complex sample surveys.

Prerequisites: The students will need to have received credit for STAT 4760/7760 or equivalent to be enrolled in this class. In other words, you will have understanding of statistical inference concepts. Having taken STAT 4310/7310 Introduction to Sampling is an advantage.

Other info: Academic integrity is fundamental to the activities and principles of a university. All members of the academic community must be confident that each person's work has been responsibly and honorably acquired, developed, and presented. Any effort to gain an advantage not given to all students is dishonest whether or not the effort is successful. The academic community regards breaches of the academic integrity rules as extremely serious matters. Sanctions for such a breach may include academic sanctions from the instructor, including failing the course for any violation, to disciplinary sanctions ranging from probation to expulsion. When in doubt about plagiarism, paraphrasing, quoting, collaboration, or any other form of cheating, consult the course instructor.

If you have special needs as addressed by the Americans with Disabilities Act (ADA) and need assistance, please notify the Office of Disability Services, A038 Brady Commons, 882-4696 or course instructor immediately. Reasonable efforts will be made to accommodate your special needs.

Grade structure: homeworks (30%) + class presentation (30%) + takehome final (40%)

The homework exercises (about 5–6 throughout the semester) will represent a mix of theoretical questions, and practical examples to be studied with the complex data sets. The class presentation (about 20–25 min) will be one of the additional topics papers, see the list below. Expect the takehome final to be all-inclusive, with theoretical and practical questions, as well as questions based on readings.

Data sets: Students on the biostat track might want to use NHANES data for their homeworks (see links below). Students from social science tracks might want to use GSS or CPS surveys. Students in education might want to use NAEP data. Other data sets might be used from the student's area of interest; those should have sufficiently complex sample design and non-trivial design effects.

Software: Design-based estimation is now incorporated in many software titles. Usability varies from the traditional set of estimators (means, totals, ratios, proportions) to multi-stage designs, and to a variety of analytical tools (linear regression, logistic regression, survival models, and other multivariate techniques). The current leaders appear to be Stata, R and (SAS-callable) SUDAAN. All of them can handle stratified clustered designs with Taylor-series linearization or jackknife and BRR replicate variance estimation, for the linear statistics and a variety of regression estimation procedures, and that is probably as far as most analytic uses of survey would go in the likely applications. A review of the existing software (although it does not seem to have been updated recently) can be found at <http://www.hcp.med.harvard.edu/statistics/survey-soft/>.

2 Content

The class will consist of several modules, as outlined below.

| | Topics | Readings |
|-------------------|--|---|
| 1. | Basic concepts: SRS, WR, WOR, stratified samples, clustered samples, (Narain-)Horwitz-Thompson estimator. Asymptotic normality | Ch. 1–3 of <i>TH97</i> , Ch. 1 and 2 of <i>SHS89</i> , Dalenius (1994), Ch. 1 of <i>LP04</i> , Ch. 2 of <i>KG99</i> , Ch. 1–2 of <i>CS05</i> ; Brewer & Donadio (2003) |
| 2. | Design-based, model-based, model-assisted, predictive approaches to survey inference | Binder & Roberts (2003), Kish & Frankel (1974), Brewer (2002), Särndal, Swensson & Wretman (1992), Ch. 3 of <i>CS05</i> |
| 3. | Survey weights | Ch. 4 of <i>KG99</i> , Sec. 6.2 of <i>TH97</i> , Pfeffermann (1993), Korn & Graubard (1995) |
| 4. | Analysis of subdomains and subpopulations | Skinner (1989) = Ch. 3 of <i>SHS89</i> , Ch. 6 of <i>LP04</i> , Bellhouse & Rao (2002), Hidiroglou & Patak (2004) |
| 5. | Nonlinear statistics, regression and estimating equations | Binder (1983), Skinner (1989), Ch. 4 and Sec. 6.4–6.5 of <i>TH97</i> , Fuller (1975), Fuller (2002), Sec. 11.2 of <i>CS05</i> |
| 6. | Missing data | Kalton & Kasprzyk (1986), Little (2003 <i>b</i>), Ch. 4 of <i>LP04</i> , Ch. 4 of <i>KG99</i> , Ch. 13 of <i>CS05</i> ; Little & Vartivarian (2005), Haziza & Rao (2006), Kim, Michael, Fuller & Kalton (2006) |
| 7. | Small area estimation | Rao (2003), Ghosh & Rao (1994), Sec. 10.4 of <i>CS05</i> ; Fay & Herriot (1979), Prasad & Rao (1990), Ghosh, Natarajan, Stroud & Carlin (1998), Lehtonen, Särndal & Veijanen (2003), special issue of <i>Statistics in Transition</i> |
| 8. | Variance estimation and resampling inference | Sec. 4.2 of <i>TH97</i> , Ch. 5 of <i>KG99</i> , Ch. 5 of <i>LP04</i> , Shao (1996), Krewski & Rao (1981), Rao & Wu (1988), Rao, Wu & Yue (1992), Ch. 7 and 9 of <i>CS05</i> |
| Additional topics | | |
| i. | Empirical likelihood inference | Chen & Qin (1993), Wu (2004), Wu & Rao (2006) |
| ii. | Multilevel models | Pfefferman, Skinner, Holmes, Goldstein & Rasbash (1998), Rabe-Hesketh & Skrondal (2006) |
| iii. | Sampling in space and time | Binder & Hidiroglou (1988), Fuller (1990), Ernst (1999), Ch. 7 of <i>TH97</i> |
| iv. | Bayesian methods | Little (2003 <i>a</i>) = Ch. 4 of <i>CS03</i> , Ghosh et al. (1998), You & Chapman (2006) |
| v. | Case-control studies | Scott & Wild (2003) = Ch. 8 of <i>CS03</i> , Ch. 9 of <i>KG99</i> |
| vi. | Disclosure risk | Skinner & Carter (2003) |
| vii. | Inverse sampling | Rao, Scott & Benhin (2003) |
| viii. | Non-sampling error | Lesser & Kalsbeek (1992) |
| ix. | Post-stratification | Holt & Smith (1979), Valliant (1993) |
| x. | Survey methodology and cognitive issues | Groves, Couper, Lepkowski, Singer & Tourangeau (2004), <i>Statistics Canada</i> (2003) |

3 Readings

The list of topics and readings should not be intimidating. This is the list of “everything-you-need-to-know-about-survey-statistics” (unless you do methodological research in the area). The readings are provided for your reference, so that you could consult your syllabus should the need arise in your practical work to get started with the literature search. The course is divided into the main part that will be delivered by the instructor, with the readings that generally

are book chapters, invited papers, or other big reviews of the topic; and the optional part, with the topics to be picked by students for their term presentation, and the readings being the research papers.

There are several great books on the topic of complex survey sampling and data analysis. Some of them, mostly earlier ones, tend to gravitate to the issues of sampling *per se* and mathematical foundations: Kish (1965), Cochran (1977), Wolter (1985), Thompson (1992), Levy & Lemeshow (2003), Chaudhuri & Stenger (2005) (referred to as *CS05* above). Other more recent books tend to focus more on the analytical methods developed to address a wide range of practical problems: Skinner, Holt & Smith (1989) [*SHS89*], Thompson (1997) [*TH97*], Korn & Graubard (1999) [*KG99*], Chambers & Skinner (2003) [*CS03*], Lehtonen & Pahkinen (2004) [*LP04*]. If you have any of those books in your library, it will cover most of the “first order” topics in the first half of the course, and some of the “second order” selective topics. A summary of historical developments in survey statistics is given in Rao (2005). There are also some highly specialized monographs, such as Särndal et al. (1992), Rao (2003) or Tillé (2006).

There is a broad range of articles published in top journals such as *Annals of Statistics*, *JASA*, *JRSSB*, *Biometrika*, but the leading journal in the field dedicated solely to survey statistics is *Survey Methodology* published by Statistics Canada.

4 Educational objectives

Upon completion of the course, the students will:

- understand the importance of design-based (randomization) inference;
- know the implications of complex sampling designs for point and interval estimation;
- by using the randomization inference paradigm, be able to compute means and variances of simple statistics;
- know and be able to verify the domains of applicability of asymptotic normality, including results for non-linear statistics;
- be aware of the subtleties that arise in variance estimation, and be able to find ways to estimate variances in difficult situations, including those with (adjustments for) non-response;
- specify the major features of complex survey designs in their favorite software;
- perform analysis of (generalized) linear models, including analysis on subdomains, with appropriate design specification;
- be aware of the broad spectrum of research problems in area of survey statistics.

5 Links

Data sets

NHANES: <http://www.cdc.gov/nchs/nhanes.htm>
GSS: <http://www.norc.umd.edu/projects/gensoc.asp>
CPS: <http://www.census.gov/cps/>
NAEP: <http://nces.ed.gov/nationsreportcard/>

Software

Stata: <http://www.stata.com/stata9/svy.html>
R: <http://cran.us.r-project.org/src/contrib/Descriptions/survey.html>
SUDAAN: <http://www.rti.org/sudaan/>

Publications

Survey Methodology journal, open access:

<http://www.statcan.ca/bsolc/english/bsolc?catno=12-001-X&CHROPG=1>

My personal set of references:

<http://www.citeulike.org/user/ctacmo/tag/survey>

ASA Survey Research Methods Section:

<http://www.amstat.org/sections/SRMS/index.html>

Statistics in Transition special issue on SAE:

<http://www.stat.gov.pl/english/sit/sit73/index.htm>

Statistics Canada MA readings:

<http://www.statcan.ca/english/employment/ma/readings.htm>

NIH references:

http://archive.nlm.nih.gov/proj/dxpnet/nhanes/docs/doc/sample_survey/references.php

References

- Bellhouse, D. R. & Rao, J. N. K. (2002), 'Analysis of domain means in complex surveys', *Journal of Statistical Planning and Inference* **102**, 47–58.
- Binder, D. A. (1983), 'On the variances of asymptotically normal estimators from complex surveys', *International Statistical Review* **51**, 279–292.
- Binder, D. A. & Hidiroglou, M. A. (1988), Sampling in time, in P. R. Krishnaiah & C. R. Rao, eds, 'Handbook of Statistics', Vol. 6, North Holland, Amsterdam, pp. 187–211.
- Binder, D. A. & Roberts, G. R. (2003), Design-based and model-based methods for estimating model parameters, in R. L. Chambers & C. J. Skinner, eds, 'Analysis of Survey Data', John Wiley & Sons, New York, chapter 3.
- Brewer, K. (2002), *Combined Survey Sampling Inference*, Arnold/Oxford University Press.
- Brewer, K. & Donadio, M. E. (2003), 'The high entropy variance of the Horvitz-Thompson estimator', *Survey Methodology* **29**(2), 189–196.
- Chambers, R. L. & Skinner, C. J., eds (2003), *Analysis of Survey Data*, Wiley series in survey methodology, Wiley, New York.
- Chaudhuri, A. & Stenger, H. (2005), *Survey Sampling: Theory and Methods*, Vol. 181 of *Statistics: Textbooks and Monographs*, 2nd edn, Chapman & Hall/CRC, Boca Raton, FL.
- Chen, J. & Qin, J. (1993), 'Empirical likelihood estimation for finite populations and the effective usage of auxiliary information', *Biometrika* **80**(1), 107–116.
- Cochran, W. G. (1977), *Sampling Techniques*, 3rd edn, John Wiley and Sons, New York.
- Dalenius, T. (1994), A first course in survey sampling, in P. R. Krishnaiah & C. R. Rao, eds, 'Handbook of Statistics: Sampling', Vol. 6, Elsevier: North Holland, chapter 2.
- Ernst, L. R. (1999), The maximization and minimization of sample overlap problems: A half century of results, Technical report, U.S. Bureau of Labor Statistics.
- Fay, R. E. & Herriot, R. A. (1979), 'Estimates of income for small places: An application of James-Stein procedures to census data', *Journal of the American Statistical Association* **74**(366), 269–277.
- Fuller, W. A. (1975), 'Regression analysis for sample survey', *Sankhya Series C* **37**, 117–132.
- Fuller, W. A. (1990), 'Analysis of repeated surveys', *Survey Methodology* **16**(2), 167–180.
- Fuller, W. A. (2002), 'Regression estimation for survey samples (with discussion)', *Survey Methodology* **28**(1), 5–23.

- Ghosh, M., Natarajan, K., Stroud, T. W. F. & Carlin, B. P. (1998), 'Generalized linear models for small-area estimation', *Journal of the American Statistical Association* **93**(441), 273–282.
- Ghosh, M. & Rao, J. N. K. (1994), 'Small area estimation: An appraisal', *Statistical Science* **9**(1), 55–76.
- Groves, R. M., Couper, M. P., Lepkowski, J. M., Singer, E. & Tourangeau, R. (2004), *Survey Methodology*, Wiley Series in Survey Methodology, John Wiley and Sons, New York.
- Haziza, D. & Rao, J. N. (2006), 'A nonresponse model approach to inference under imputation for missing survey data', *Survey Methodology* **32**(1), 53–64.
- Hidiroglou, M. A. & Patak, Z. (2004), 'Domain estimation using linear regression', *Survey Methodology* **30**, 67–78.
- Holt, D. & Smith, T. M. F. (1979), 'Post stratification', *Journal of the Royal Statistical Society, Series A* **142**(1), 33–46.
- Kalton, G. & Kasprzyk, D. (1986), 'The treatment of missing survey data', *Survey Methodology* **12**(1), 1–16.
- Kim, J. K., Michael, B. J., Fuller, W. A. & Kalton, G. (2006), 'On the bias of the multiple-imputation variance estimator in survey sampling', *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **68**(3), 509–521.
- Kish, L. (1965), *Survey Sampling*, John Wiley and Sons, New York.
- Kish, L. & Frankel, M. R. (1974), 'Inference from complex samples', *Journal of the Royal Statistical Society, Series B* **36**, 1–37.
- Korn, E. L. & Graubard, B. I. (1995), 'Analysis of large health surveys: Accounting for the sampling design', *Journal of the Royal Statistical Society, Series A* **158**(2), 263–295.
- Korn, E. L. & Graubard, B. I. (1999), *Analysis of Health Surveys*, John Wiley and Sons.
- Krewski, D. & Rao, J. N. K. (1981), 'Inference from stratified samples: Properties of the linearization, jackknife and balanced repeated replication methods', *The Annals of Statistics* **9**(5), 1010–1019.
- Lehtonen, R. & Pahkinen, E. (2004), *Practical Methods for Design and Analysis of Complex Surveys*, Statistics in Practice, 2nd edn, John Wiley & Sons, New York.
- Lehtonen, R., Särndal, C.-E. & Veijanen, A. (2003), 'The effect of model choice in estimation for domains, including small domains', *Survey Methodology* **29**(1), 33–44.
- Lesser, V. M. & Kalsbeek, W. D. (1992), *Non-sampling Error in Surveys*, John Wiley and Sons, New York.
- Levy, P. S. & Lemeshow, S. (2003), *Sampling of Populations: Methods and Applications*, 3rd edn, John Wiley & Sons, New York.
- Little, R. J. (2003a), The Bayesian approach to sample survey inference, in R. Chambers & C. J. Skinner, eds, 'Analysis of Survey Data', John Wiley and Sons, chapter 4.
- Little, R. J. (2003b), Bayesian methods for unit and item nonresponse, in R. Chambers & C. J. Skinner, eds, 'Analysis of Survey Data', John Wiley and Sons, chapter 18.
- Little, R. J. & Vartivarian, S. (2005), 'Does weighting for nonresponse increase the variance of survey means?', *Survey Methodology* **31**(2), 161–168.
- Pfefferman, D., Skinner, C. J., Holmes, D. J., Goldstein, H. & Rasbash, J. (1998), 'Weighting for unequal selection probabilities in multilevel models', *Journal of Royal Statistical Society, Series B* **60**(1), 23–40.
- Pfeffermann, D. (1993), 'The role of sampling weights when modeling survey data', *International Statistical Review* **61**, 317–337.

- Prasad, N. G. N. & Rao, J. N. K. (1990), 'The estimation of the mean squared error of small-area estimators', *Journal of the American Statistical Association* **85**(409), 163–171.
- Rabe-Hesketh, S. & Skrondal, A. (2006), 'Multilevel modelling of complex survey data', *Journal of the Royal Statistical Society, Series A* **169**(4).
- Rao, J. N. K. (2003), *Small Area Estimation*, Wiley series in survey methodology, John Wiley and Sons, New York.
- Rao, J. N. K. (2005), 'Interplay between sample survey theory and practice: An appraisal', *Survey Methodology* **31**(2), 117–138.
- Rao, J. N. K. & Wu, C. F. J. (1988), 'Resampling inference with complex survey data', *Journal of the American Statistical Association* **83**(401), 231–241.
- Rao, J. N. K., Wu, C. F. J. & Yue, K. (1992), 'Some recent work on resampling methods for complex surveys', *Survey Methodology* **18**(2), 209–217.
- Rao, J., Scott, A. & Benhin, E. (2003), 'Undoing complex survey data structures: Some theory and applications of inverse sampling (with discussion)', *Survey Methodology* **29**(2), 107–128.
- Särndal, C.-E., Swensson, B. & Wretman, J. (1992), *Model Assisted Survey Sampling*, Springer, New York.
- Scott, A. & Wild, C. (2003), Fitting logistic regression models in case-control studies with complex sampling, in R. Chambers & C. J. Skinner, eds, 'Analysis of Survey Data', John Wiley and Sons, chapter 8.
- Shao, J. (1996), 'Resampling methods in sample surveys', *Statistics* **27**, 203–254. with discussion.
- Skinner, C. & Carter, R. (2003), 'Estimation of a measure of disclosure risk for survey microdata under unequal probability sampling', *Survey Methodology* **29**(2), 177–180.
- Skinner, C. J. (1989), Domain means, regression and multivariate analysis, in C. J. Skinner, D. Holt & T. M. Smith, eds, 'Analysis of Complex Surveys', Wiley, New York, chapter 3, pp. 59–88.
- Skinner, C. J., Holt, D. & Smith, T. M., eds (1989), *Analysis of Complex Surveys*, Wiley, New York.
- Statistics Canada (2003), *Survey Methods and Practices*, Ottawa. Catalogue No. 12-587-XPE.
- Thompson, M. E. (1997), *Theory of Sample Surveys*, Vol. 74 of *Monographs on Statistics and Applied Probability*, Chapman & Hall/CRC, New York.
- Thompson, S. K. (1992), *Sampling*, John Wiley and Sons, New York.
- Tillé, Y. (2006), *Sampling Algorithms*, Springer Series in Statistics, Springer, New York.
- Valliant, R. (1993), 'Poststratification and conditional variance estimation', *Journal of the American Statistical Association* **88**(421), 89–96.
- Wolter, K. M. (1985), *Introduction to Variance Estimation*, Springer, New York.
- Wu, C. (2004), 'Some algorithmic aspects of the empirical likelihood method in survey sampling', *Statistica Sinica* **14**, 1057–1067.
- Wu, C. & Rao, J. N. K. (2006), 'Pseudo-empirical likelihood ratio confidence intervals for complex surveys', *The Canadian Journal of Statistics/La revue canadienne de statistique* **34**(3), 359–376.
- You, Y. & Chapman, B. (2006), 'Small area estimation using area level models and estimated sampling variances', *Survey Methodology* **32**(1), 97–103.

The references are also available at <http://www.citeulike.org/user/ctacmo/tag/stat9100svy>

STAT 9100.3 COMPLEX SURVEYS - INTRO

Note Title

1/11/2007

- INTRODUCE MYSELF
- HAND OUT SYLLABUS
- MECHANICS:
 - CLASS PERIODS
 - STUDENT TALKS
 - HOME ASSIGNMENTS
 - TAKEHOME FINAL
- UNIQUENESS OF THE COURSE!
 - + SHORTAGE OF SURVEY STATISTICIANS
 - + WIDE USE OF SURVEY STATS
 - FEDERAL: BLS, CENSUS; STAT CANADA
 - HEALTH: CDC, NCI, ... - HEALTH DATA
 - PRIVATE: RTI, NORC, WESTAT, ...
 - GOD ONLY KNOWS WHAT ELSE!
 - + MY OWN RESEARCH INTERESTS
- TOPICS IN THE SYLLABUS :
 - BASIC TOPICS IN CLASS BY ME
 - ADVANCED TOPICS: SELECT PAPERS
 - SHOW OF HANDS - SUGGEST TOPICS
- BROAD OBJECTIVES:
 - + OVERALL FAMILIARITY WITH WHAT SURV STAT IS ABOUT
 - + UNDERSTANDING OF THE DESIGN-BASED ESTIM
- TO THAT EFFECT, A HUGE LIST OF REFS,
CLICKABLE LINKS IN SYLLABUS

BOOKS ON SAMPLING:

- THOMPSON, M (1997) - THE MAIN TEXT
INTERMEDIATE TO UPPER LEVEL,
QUITE RIGHT FOR US
- LEHTONEN & PARKINEN (2004)
KORN & GRAUBARD (1999)
 - SOMEWHAT LIGHTER BUT STILL RIGOROUS
- KISH, COCHRAN, THOMPSON (92), LEVY & LEMESHOW
 - GOOD OL' BOOKS, DON'T HAVE
NEWER STUFF, RATHER DRY
- SHS, CHAMBERS & SKINNER -
COLLECTIONS OF RESEARCH MONOGRAPHS,
THE MOST ADVANCED & UP TO DATE
- SÄRNDAL, SWENSSON, WRETMAN -
MORE DIFFICULT TO DIGEST
- SPECIALIZED MONOGRAPHS: RAO; TILLÉ; ...

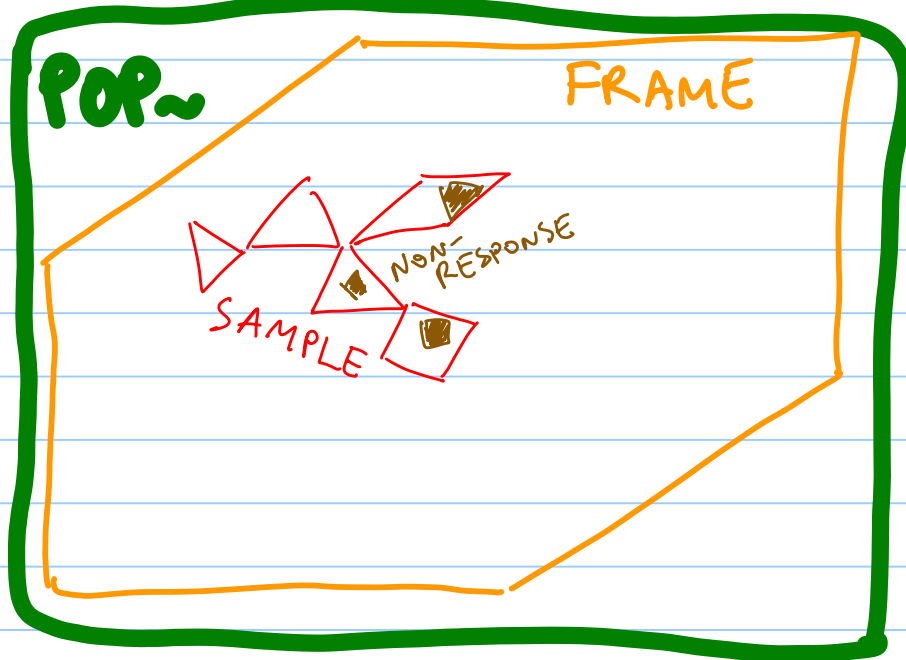
FIG. FROM CH 1 OF LP(04) -
SURVEY PROCESS.

WHERE DO...

- STATISTICIANS
- COGNITIVE / PSYCHOMETRIC PEOPLE
- MANAGERS
- FIELD STAFF
- CLIENTS

... FIT IN?

SAMPLING PROCESS:



NOTE THAT POP~S TO WHICH THE RESEARCHER WANTS TO GENERALIZE MAY NOT BE THOSE FOR WHICH (LARGE SCALE) DATA WERE COLLECTED !
TARGET VS. REPRESENTED POP~

OBSERVED SUMMARY $Y =$

$=$ TRUE Y +

+ $(Y[\text{POP~}] - Y[\text{FRAME}])$ +

+ RANDOM~ ERROR

+ INSTRUMENT~ ERROR

+ NON-RESPONSE ERROR

+ ADJUSTMENTS

UNDERCOVERAGE

SAMPLING ERROR

MEASUREMENT

MEASUREMENT

DESIGN PARADIGM :

THE RANDOM COMPONENT IS DUE TO RANDOM~, BUT IT IS NOT THE CHARACTERISTIC ITSELF.

DESCRIPTIVE USE : TOTALS, MEANS, RATIOS, PROPORTIONS

ANALYTICAL USE : STAT~ REL~ MODELS

(LINEAR, LOGISTIC, SVY REGRESSION;

MULTILEVEL & STRUCTURAL EQ~ MODELS)

EXAMPLE OF A REAL COMPLEX SVT:
NHANES III (K99)

RUN INTRO-NHANES.DO

<http://www.cdc.gov/nchs/nhanes.htm>

NHANES III :

2812 PSUs

↳ 13 WP1

↳ 2799

↳ 34 STRATA, ALLOW OVERSAMPLING
OF NON-WHITES

↳ 2 PSU/STRATUM

I

89 STANDS FOR MOBILE EXAM CENTERS

↳ SEGMENTS

II

↳ LIST HHs IN EACH SEGMENT

IV

↳ INDIVIDUALS \subseteq HH

SCREENING, DEMOGRAPHY

⇒ ≈ 1 IN 5 SAMPLED HH

CONTRIBUTED OBSERVATION UNITS

INSTRUMENTS:

- INDIVIDUAL Q'RE

- BLOOD PRESSURE MEAS: HH

- FAMILY Q'RE

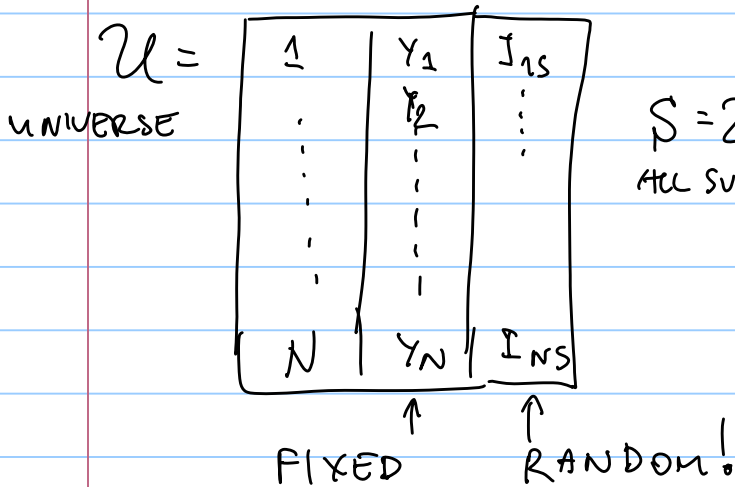
- MEC → AUTOMATED DATA COLLECTION

STAT 9100 SVY - SR, S & SIMPLE DESIGNS

Note Title

1/12/2007

FRAMEWORK & DEFINITIONS



$S = 2^U$
ALL SUBSETS

CAPITALS / LOWERCASE
POP~ / SAMPLE

$$T[Y] = \sum_{j=1}^N Y_j \quad \text{--- POP~ TOTAL}$$

$s = (j_1, \dots, j_n)$ - SAMPLE (UNORDERED - SUFF STAT)

PROBABILITY SPACE / SAMPLE DESIGN: $p: S \rightarrow [0,1]$
 $\forall s \in S$ $p(s)$ IS PROB. OF SELECTING THIS PARTICULAR SAMPLE

EXAMPLES: $U = \{1, 2, 3\}$ SRS WR, SRSWOR

SAMPLE SIZE: $n(s) = \sum_{j=1}^N I_{js} = \# \text{ UNITS IN } s$

FIXED SIZE DESIGNS: $p(s) = 0 \iff n(s) \neq n$

INCLUSION PROBITY: $\pi_j = E_{p(s)} I_{js}$

NOTE THAT $\sum_j \pi_j = E n(s) > 1$

SELF - WGT DESIGN / EPSEM: $\forall j \quad \pi_j = \text{const} = n/N$
 NOT EVERY EPSEM IS NICE!

JOINT INCLUSION PROBITY: $T_{ijk} = \sum_{p(s)} I_{js} I_{ks}$
 IMPORTANT FOR VARIANCE ESTIM~!

SRSWOR:

$$\forall s \quad n(s) = n, \quad p(s) = \text{const} = 1 / \binom{N}{n}$$

$\forall j$, # SAMPLES WHERE IT IS USED = $\binom{N-1}{n-1}$
 HENCE, $\bar{j}_j = \frac{\binom{N-1}{n-1}}{\binom{N}{n}} = \frac{(N-1)!}{(n-1)!(N-n)!} \cdot \frac{n!(N-n)!}{N!} = \frac{n}{N}$ - ESEM

JOINT INCLUSION PROBITY:

\forall PAIR j, k # SAMPLES = $\binom{N-2}{n-2}$

HENCE, $\bar{j}_{jk} = \frac{\binom{N-2}{n-2}}{\binom{N}{n}} = \frac{(N-2)!}{(n-2)!(N-n)!} \cdot \frac{n!(N-n)!}{N!} = \frac{n(n-1)}{N(N-1)}$

\bar{j}_{jk} vs n^2/N^2 ?

$$\frac{n(n-1)}{N(N-1)} < \frac{n^2}{N^2} \Leftrightarrow \frac{n-1}{N-1} < \frac{n}{N} \Leftrightarrow nN - N < nN - n$$

$n < N$

FINITE POPULATION CONSISTENCY:

THE ESTIMATOR = THE POP~ PARAM
 WHEN $s = \mathcal{U}$

EXPECTATIONS

$$\begin{aligned} E_{p(s)} \sum_{j \in S} z_j &= \sum_{s \in S} p(s) \sum_{j \in S} z_j = \\ &= \sum_{s \in S} p(s) \sum_{j=1}^N z_j I_{js} = \sum_{j=1}^N \left\{ z_j \sum_{s \in S} p(s) I_{js} \right\} = \sum_{j=1}^N z_j \pi_j \end{aligned}$$

$$z_j = 1 \Rightarrow \text{SAMPLE SIZE} \quad E n(s) = \sum_{j=1}^N \pi_j$$

$$\left. \begin{array}{l} z_j \text{ IS SOMETHING REAL} \\ \& \\ \text{EPSEM} \end{array} \right\} \Rightarrow E \sum_{j \in S} y_j = \sum_{j=1}^N y_j \frac{\pi_j}{N} = \frac{1}{N} T[y]$$

$$\bar{y}_n = \frac{1}{n} \sum_{j \in S} y_j \Rightarrow E \bar{y} = T[y]/N - \text{UNBIASED!}$$

EXPANSION ESTIMATOR OF TOTAL:

$$t[y] = N \bar{y}_n \quad \text{FOR EPSEM DESIGNS}$$

FOR ARBITRARY DESIGNS,
HORVITZ-THOMPSON (-NARAIN) ESTIMATOR

$$t_{HT}[y] = \sum_{j \in S} y_j / \pi_j$$

\Rightarrow SHOW IT IS UNBIASED — IN CLASS EXERCISE?
— SHORT WORK PROBLEM?

| | $p(s)$ | π_j | π_{jk} | t_{HT} | $\forall t_{HT}$ |
|-------------|--------|---------|------------|----------|------------------|
| $\{1\}$ | 0 | | | | |
| $\{2\}$ | 0 | | | | |
| $\{3\}$ | 0 | | | | |
| $\{1,2\}$ | 3/12 | | | | COMPARE TO SRS |
| $\{1,3\}$ | 4/12 | | | | |
| $\{2,3\}$ | 5/12 | | | | |
| $\{1,2,3\}$ | 0 | | | | |

DESIGN VARIANCE :

$$\begin{aligned}
 V\left(\sum_{j \in S} z_j\right) &= V_{p(s)}\left(\sum_{j=1}^N z_j I_{js}\right) \\
 &= \sum_{j=1}^N z_j^2 V I_{js} + \sum_{j \neq k} z_j z_k \text{Cov}(I_{js}, I_{ks}) = \\
 &= \sum_{j=1}^N z_j^2 \pi_j (1 - \pi_j) + \sum_{j \neq k} z_j z_k (\pi_{jk} - \pi_j \pi_k)
 \end{aligned}$$

$$\text{Cov}(I_{js}, I_{ks}) = E I_{js} I_{ks} - E I_{js} E I_{ks} = \pi_{jk} - \pi_j \pi_k$$

FIXED SIZE: $P(\sum I_{js} = n) = 1$

$$\begin{aligned}
 \sum_{k \neq j} \text{Cov}(I_{js}, I_{ks}) &= \text{Cov}(I_{js}, n - I_{js}) = -V(I_{js}) \\
 \Rightarrow V\left(\sum_{j \in S} z_j\right) &= \frac{1}{2} \sum_{j \neq k} (z_j - z_k)^2 (\pi_j \pi_k - \pi_{jk})
 \end{aligned}$$

$$\text{HT: } z_j = y_j / \pi_j \Rightarrow V(t_{\text{HT}}(y)) = \sum_{j=1}^N y_j^2 \left(\frac{1}{\pi_j} - 1\right) + \sum_{j \neq k} y_j y_k \left(\frac{\pi_{jk}}{\pi_j \pi_k} - 1\right)$$

$$\Rightarrow \text{FIXED } n = \frac{1}{2} \sum_{j \neq k} \left(\frac{y_j}{\pi_j} - \frac{y_k}{\pi_k}\right)^2 \Omega_{jk}, \quad \Omega_{jk} = \pi_j \pi_k - \pi_{jk}$$

? → DERIVE THE VARIANCE FORMULA FOR SRS
 ⇒ DISCUSS THE FINITE POP~ CORR~
 (MENTION V-STAT PRES~ OF V?)

$$\Omega_{jk} = n^2/N^2 - n(n-1)/N(N-1) = \frac{n^2(N-1) - Nn(n-1)}{N^2(N-1)} = \frac{n(N-n)}{N^2(N-1)}$$

DESIGN EFFECT:

$$\text{DEFF}(t) = \frac{V_{p(s)}[t]}{V_{\text{SRS}}[t]}$$

↑
WOR

ESTIMATORS OF VARIANCE

IF z_{jk} IS A CHARACTERISTIC OF A PAIR OF OBSERVATIONS, THEN AN ESTIMATE OF $T(z_{jk})$

$$t(z) = \sum_{j \neq k} z_{jk} / \pi_{jk}$$

\Rightarrow

$$v(t_{HT}) = \frac{1}{2} \sum_{\substack{j \neq k \\ j \neq k}} \left(\frac{y_j}{n_j} - \frac{y_k}{n_k} \right)^2 \frac{\pi_j \pi_k - \pi_{jk}}{\pi_{jk}}$$

YATES-GRUNDY-SEN

ESTIMATOR

HT-SRVE
ESTIMATOR

FOR A U-STAT
OF ORDER 2

UNBIASED IF $\forall j, k \pi_{jk} > 0$

WEIRD IF $\pi_{jk} = 0!$

OTHER TRADITIONAL DESIGNS

SRS WITH REPLACEMENT

↳ n UNITS WITH REPLACEMENT

$$1 - \pi_j = \frac{(N-1)^n}{N^n} \sim \frac{N^n - nN^{n-1} + \frac{n(n-1)}{2}N^{n-2} - \dots}{N^n}$$

$$\pi_j \approx n/N + o(n/N) \quad \text{as } n/N \rightarrow 0$$

$$1 - \pi_{ij} = (N-2)^n / N^n = 1 - \frac{2n}{N} + \frac{2n(n-1)}{N^2} - \dots$$

NICE DESIGN: SRSWR + SRSWR = SRSWR!

NOT TRUE FOR WOR

$$E u(s) = \sum_j \pi_j = n - \frac{n(n-1)}{2N} + \dots$$

$$V [t_{HT}(y)] = \sum_{j=1}^N y_j^2 \left(\frac{1-\pi_j}{\pi_j} \right) + \sum_{j \neq k} y_j y_k \left(\frac{\pi_{jk} - \pi_j \pi_k}{\pi_j \pi_k} \right) = \dots$$

CLUSTERED SAMPLING

$U = B_1 \cup \dots \cup B_L$, L CLUSTERS ARE SELECTED AT RANDOM, EVERY UNIT OBSERVED

$$\pi_j = \pi_{jk} = l/L, \quad j, k \in B_t \quad \rho_{jk} = l/L \left(\frac{l}{L} - 1 \right) < 0$$

$$\pi_{jk} = \frac{l(l-1)}{L(L-1)}, \quad j \in B_t, k \in B_{t'}, t \neq t' \quad \rho_{jk} = \frac{l(l-l)}{L^2(L-1)} > 0$$

WHAT IS AN EFFICIENT STRATEGY?

STRATIFIED SAMPLING

$$\mathcal{U} = \mathcal{S}_1 \cup \dots \cup \mathcal{S}_H, \quad |\mathcal{S}_k| = N_k \text{ known}$$

WITHIN EACH \mathcal{S}_h , AN SRS IS TAKEN:

$$\pi_j = n_h / N_h \quad \text{FOR } j \in \mathcal{S}_h$$

$$t_{HT}(Y) = \sum_h^H N_h \bar{Y}_h; \quad \bar{Y}_{STR} = \sum_{h=1}^H W_h \bar{Y}_h$$

$$\bar{Y}_h = \frac{1}{n_k} \sum_{\substack{j \in \mathcal{S}_h \\ j \in \mathcal{S}}} Y_j \quad W_h = N_h / \sum_k N_k$$

$$V[t_{HT}] = ?$$

$$\rho_{jk} = 0, \quad j \in \mathcal{S}_h, k \in \mathcal{S}_{h'}, h \neq h'$$

$$\rho_{jk} = \text{TRUE FOR SRS, } j, k \in \mathcal{S}_h$$

(2.46)

(2.47)

\Rightarrow WHAT IS THE OPTIMAL STRATIF~ STRATEGY?

hw: Neyman Alloc~

$$\forall j \in \mathcal{S}_h \text{ COST} = C_h$$

$$\text{TOTAL BUDGET} = C$$

$$\text{OPTIMAL ALLOCATION} = ?$$

MY PAPER:
BUDGET-OPTIMAL
ALLOC~ FOR
REPEATED
CLUSTERED SUR

OTHER ESTIMATORS OF THE TOTAL

LINEAR ESTIMATORS $e = \sum_{j \in S} d_{js} Y_j$, $\sum_{s: j \in S} p(s) d_{js} = 1$
 $= \sum_{j=1}^N d_{js} I_{js} Y_j$; $\sum_s p(s) I_{js} d_{js} = 1$

$$V_{p(s)}[e] = \sum_{j=1}^N a_j Y_j^2 + \sum_{j \neq k} a_{jk} Y_j Y_k$$

$$a_j = V[d_{js} I_{js}] = \sum_s p(s) d_{js}^2 I_{js}^{-1}$$
$$a_{jk} = \text{Cov}(d_{js} I_{js}, d_{ks} I_{ks})$$

$$v(e) = \sum_{j \in S} \frac{a_j}{\pi_j} Y_j^2 + \sum_{j \neq k} \frac{a_{jk}}{\pi_{jk}} Y_j Y_k$$

CALIBR~ FOR ARRAY w :

$$\forall \text{ a.e. } s \quad \sum d_{js} w_j = T[w]$$

$$\Rightarrow \text{MSE}(e) = -\frac{1}{2} \sum_{j \neq k} a_{jk} \left(\frac{Y_j}{w_j} - \frac{Y_k}{w_k} \right)^2 w_j w_k$$

$$a_{jk} = E[(d_{js} I_{js} - 1)(d_{ks} I_{ks} - 1)]$$

$$E_{p(s)} e = T(y) \Rightarrow$$

$$v(e) = -\frac{1}{2} \sum_{j \neq k} \left(d_{js} d_{ks} - \frac{1}{\pi_{jk}} \right) \left(\frac{Y_j}{w_j} - \frac{Y_k}{w_k} \right)^2 w_j w_k$$

PPS SAMPLING:

$$\pi_j \sim Y_j \Rightarrow \text{MSE}(e) = 0$$

STAT 9100 SVY - MULTISTAGE

Note Title

1/25/2007

1ST STAGE:

RANDOM SELECTION OF
PRIMARY SAMPLING UNITS (PSUs)
(TYPICALLY, TIPS)

↳ 2ND STAGE:

RANDOM SELECTION OF SSUs

:

OBSERVATION UNIT: THE ONE
THAT IS CONTACTED, AND
RESPONDS TO THE SVY

BIRTHDAY METHOD FOR WITHIN HH SAMPLING

$$U = \mathcal{B}_1 U \dots U \mathcal{B}_L$$
$$|\mathcal{B}_r| = M_r, \quad \sum_{r=1}^L M_r = N$$

1ST STAGE:

SAMPLE S_B OF PSU LABELS

2ND STAGE: $\forall r \in S_B,$

SAMPLE S_r IS TAKEN FROM $\mathcal{B}_r,$

$$|S_r| = m(S_r)$$

$$s = \bigcup_{r \in S_B} S_r, \quad n(s) = \sum_{r \in S_B} m(S_r)$$

1ST STAGE INCLUSION PROBABILITY:
 $P[r \in S_B] = \pi_r$

2ND STAGE INCLUSION PROBABILITY:
 $P[j \in S_r | r \in S_B] = \pi_{j|r}$

UNIT INCL. PROBABILITY:
 $\pi_j = \pi_r \pi_{j|r}$

HAT ESTIMATOR
 $t(y) = \sum_r \frac{t_r(y)}{\pi_r}$
 $t_r(y) = \sum_{j \in S_r} \pi_{j|r} y_j$
 HAT ESTIMATOR
 OF THE r-th PSU TOTAL

EXAMPLE:

1) $\pi_r = M_r / N \leftarrow \pi_{PS}$
 2) $\pi_{j|r} = m_r / M_r \leftarrow \pi_{SR}$ } $\pi_j = \pi_r \pi_{j|r} = \frac{m_r}{N}$

EPSEM: $m_r = m \forall r$
 $t(y) = \sum_{r \in S_B} \frac{N}{M_r} \sum_{j \in S_r} \frac{y_j M_r}{m} = \frac{N}{n} \sum_r \sum_j y_j = N \bar{y} = N \cdot \frac{1}{L} \sum_r \bar{y}_r$

MORE COMPLEX

ESTIMATOR $e = \sum_{r \in S_B} d_r(s_B) t_r$
 (CALIBRATION WEIGHTS) ESTIMATOR OF TOTAL IN rTH PSU

$E[e] = E_I \left[E_{II}(e | s_B) \right]$
 1ST STAGE DESIGN 2ND STAGE DESIGN

$E_{II}(t_r | s_B) = T_r \Rightarrow E[e] = E_I \sum_r d_r(s_B) T_r = \sum_{r=1}^L T_r = T(y)$
 USUALLY WANT THIS TO HOLD

$$V[e] = E_{\underline{I}} V_{\underline{\Pi}}(e | s_B) + V_{\underline{I}} E_{\underline{\Pi}}(e | s_B) =$$

$$= E_{\underline{I}} \left\{ \sum_{r \in s_B} d_r^2(s_B) V_{\underline{\Pi}}(t_r | s_B) \right\} +$$

$$+ V_{\underline{I}} \left(\sum_r d_r(s_B) T_r \right)$$

E.G. $\underline{I} = nps$, $\underline{\Pi} = SRS$:

$$\Rightarrow t[y] = N \bar{y}$$

$$V[N \bar{y}] = \frac{N^2}{e^2} \sum_r \frac{\pi_r}{m_r} \left(1 - \frac{m_r}{M_r}\right) S_r^2$$

$$+ \frac{1}{2} \sum_{r \neq q} (\pi_r \pi_q - \pi_{rq}) \left(\frac{T_r}{\pi_r} - \frac{T_q}{\pi_q} \right)^2$$

V ESTIM

NEED TO BUILD UP $j \rightarrow r \rightarrow S$

HENCE A BACKWARDS DECOMPOSITION MIGHT BE USEFUL:

$$V[e] = E_{\underline{\Pi}} V_{\underline{I}}(e | s_r, r=1, \dots, L) + V_{\underline{\Pi}} E_{\underline{I}}(e | s_r, r=1, \dots, L)$$

\nwarrow ESSENTIALLY $V_{\underline{\Pi}}[t_r]$

$$d_r(s_B) = 1/\pi_r, \quad l = \text{FIXED}$$

$$\Rightarrow \text{1ST TERM} = \frac{1}{2} \sum_{r \neq q} (\pi_r \pi_q - \pi_{rq}) \left(\frac{t_r}{\pi_r} - \frac{t_q}{\pi_q} \right)^2$$

HENCE GENERALLY

$$V[e] = V_{\underline{I}} + \sum_{r \in s_B} d_r^*(s_B) V_{\underline{\Pi}, r}$$

$\underline{I} = nps$, $\underline{\Pi} = SRS$ OF SIZE m_r

\Rightarrow

$$V_{\underline{I}} = \frac{1}{2} \sum_{r \neq q \in s_B} \left(\frac{\pi_r \pi_q - \pi_{rq}}{\pi_{rq}} \right) \left(\frac{t_r}{\pi_r} - \frac{t_q}{\pi_q} \right)^2$$

$$d_r^*(s_B) = 1/\pi_r, \quad V_{\underline{\Pi}, r} = \frac{M_r^2}{m_r} \left(1 - \frac{m_r}{M_r}\right) \frac{1}{m_r - 1} \sum_{j \in s_r} (y_j - \bar{y}_r)^2$$

* LCL, HIGH ENTROPY WITH SELECTION PROBABILITY π_r/e

$$\Rightarrow \pi_{rq} \approx \pi_r \pi_q (e-1)/e$$

$$(\pi_r \pi_q - \pi_{rq}) / \pi_{rq} \sim 1/(e-1)$$

$$\Rightarrow \sigma[t(y)] \approx \frac{1}{2} \frac{1}{e-1} \sum_{r \neq q} \left(\frac{t_r}{\pi_r} - \frac{t_q}{\pi_q} \right)^2 \sim \frac{N^2}{e} \left[\frac{1}{e-1} \sum_r (\bar{y}_r - \bar{y})^2 \right]$$

ONLY DEPENDS ON $\bar{y}_r!$

* LPOY:

CLUSTERS OF EQ SIZES $M \rightarrow m \forall r$

$$\sigma_{cl}[t] = (ML)^2 \left[\left(1 - \frac{e}{L}\right) \frac{S_b^2}{e} + \left(1 - \frac{m}{M}\right) \frac{S_w^2}{me} \right]$$

$$S_b^2 = \frac{1}{L-1} \sum_{r=1}^L (\bar{y}_r - \bar{y})^2; \quad S_w^2 = \frac{1}{L(M-1)} \sum_r \sum_j (y_j - \bar{y}_r)^2$$

$$S_b^2 = \frac{1}{e-1} \sum_{r=1}^e (y_r - \bar{y})^2; \quad S_w^2 = \frac{1}{e(m-1)} \sum_r \sum_j (y_j - \bar{y}_r)^2$$

$$\sigma = (ML)^2 \left[\left(1 - \frac{e}{L}\right) \frac{S_b^2}{e} + \left(1 - \frac{m}{M}\right) \frac{e}{L} \frac{S_w^2}{em} \right]$$

LCL $\Rightarrow 2^{ND}$ PERM IS SMALL

STAT 9100 SVY - PPS SAMPLING / ESTIM

Note Title

1/25/2007

JTPS - PROB PROPORTIONAL TO SIZE

WHY? LOOK AT $(y_j/n_j - y_k/n_k)^2$ TERMS!

MORE GENERAL: SELECT PSU'S WITH PROB
 $\pi_r = l \alpha_r$, $\sum \alpha_r = 1$, $\alpha_r \leq 1/e$ (OTHERWISE INCLUDE w_{p1})

SAMPLING

OFTEN THE CASE IS THAT $n=2$, SO
MOST SAMPLING SCHEMES WORK EASILY
WITH $n=2$, WITH POSSIBLE EXTENSIONS TO $n>2$

ASSUME WE CAN SELECT ONE UNIT WITH
GIVEN PROBTY.

THE JOINT PROBABIES OF SELECTION, AND
HENCE $\psi[t]$ VARY FROM ONE SCHEME TO ANOTHER!

(*)

SUCCESSIVE SAMPLING

DRAW l TIMES W/O REPLACEMENT:

1ST DRAW: UNIT r , PROB p_r

2ND DRAW: UNIT q , PROB $p_q / (1 - p_r)$

$$\pi_r = p_r \left(1 + \sum_{q \neq r} \frac{p_q}{1 - p_q} \right), \pi_{rq} = p_r p_q \left(\frac{1}{1 - p_r} + \frac{1}{1 - p_q} \right)$$

FIND p_r : $\pi_r = l \alpha_r \leftarrow$ ITERATIVE PROCEDURE?

(*)

FELLEGI'S METHOD

(*)

1ST DRAW: $p_r = d_r$

2ND DRAW: $\frac{p_g}{1-p_r}$ WHERE $\forall g \sum_{t \neq g} \frac{\alpha_t p_g}{1-p_t} = \alpha_g$

ROTATING DESIGNS: REPEAT
THE 2ND DRAW STRATEGY

(*)

SAMPFORD'S METHOD

1ST DRAW: r , $p_r = d_r$

2ND, ..., LTH DRAW: WITH REPLACEMENT, $p_r \propto \frac{d_r}{1-d_r}$

REJECT & RESTART IF REPEATED UNITS

$e \uparrow \Rightarrow$ REJECT MORE OFTEN, TAKES LONGER-
IMPORTANT FOR SIMUL!

MADDOX'S ORDERED SYSTEMATIC PROCEDURE

$$[0, 1] = [0, d_1] \cup [d_1, d_1 + d_2] \cup \dots \cup [1 - d_L, 1]$$

TAKE $U \sim U(0, 1/l)$

UNIT r IS SELECTED $\Leftrightarrow \sum_{j=0}^{r-1} d_j < U + \frac{1}{l} \leq \sum_{j=0}^r d_j, \xi = 0, 1, \dots, l$

AS IS THE CASE WITH SYSTEMATIC SAMPLING,
SOME $\pi_{rq} = 0 \Rightarrow$ UNBIASED IS NOT FEASIBLE

RANDOM SYSTEMATIC PROCEDURE

RANDOMIZE THE ORDER OF UNITS FIRST!

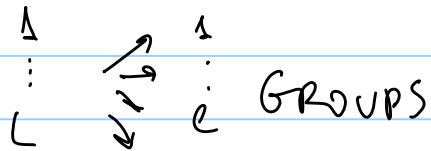
$\pi_{rq} =$ DIFFICULT \therefore

SAMPLING WITH REPLACEMENT

SAMPLE WR @ 1ST STAGE

SAMPLE II @ SUBSER STAGES

RAO-HARTLEY-COCHRAN



"TOTAL MASS" OF THE j TH GROUP : $P_j = \sum_{r \in j \text{th GROUP}} \alpha_j$

SELECT 1 PSU $r(j)$ FROM EACH GROUP $\Rightarrow \alpha_r / P_j$

CONDV HT ESTIMATOR: $e = \sum_j t_{r(j)} P_j / \alpha_{r(j)}$

APPROX VARIANCE ESTIMATORS AVAILABLE

BERNOULLI SAMPLING

$V = 1, \dots, L$ TAKE $I_{rs} \sim \text{Bernoulli}(e_{sr})$

DISADVANTAGE: RANDOM $n(s)$

SIMPLE REJECTIVE SAMPLING

- SELECT e UNITS WR WITH SELECTION PROB p_r ,
OR

- TAKE BERNOULLI SAMPLE WITH PROBAB λ_r
REJECT IF $n(s) \neq e$

NBE ENTROPY PROPERTIES.

APPROX~ OF VARIANCE

OF HT ESTIMATOR

$$\pi_{rq} \approx \pi_r \pi_q \frac{e-1}{e} \quad \text{FOR (*) DESIGNS}$$

$$\text{SRSWOR: } \pi_{jk} = \pi_j \pi_k \frac{N(n-1)}{n(N-1)} \approx \pi_j \pi_k$$

$$\Rightarrow \Delta_{jk} \approx \pi_j \pi_k \frac{N-n}{n(N-1)} \sim O(nN^{-2})$$

BREWER & DONADIO (2003)

$$\pi_{jk} \approx \pi_j \pi_k \frac{c_j + c_k}{2}$$

$c_j \leftarrow \#$ OF CHOICES; SLIGHTLY BETTER c_j :
(18) OF B & D (2003)

EFFICIENCY OF PPS SAMPLING (LPO4):

$$V_{\text{SRS}}[t] - V_{\text{PPS}}[t] = \frac{N^2}{n} G_{\alpha}(\alpha, \gamma^2/\alpha)$$

STAT 9100 SVY : IMPORTANT STAT RESULTS

Note Title

1/20/2007

CHEBYSHEV'S INEQUALITY

$$\text{MSE}_{p(s)} [t(Y)] = E_{p(s)} (t(Y) - T(Y))^2 =$$

$$= \text{FOR SOME } c, \sum_{s: |t(Y) - T(Y)| < c} ()^2 + \sum_{s: |t(Y) - T(Y)| \geq c} ()^2 \geq$$

$$\geq c^2 P \{ |t(Y) - T(Y)| \geq c \}$$

$$\Rightarrow P \{ |t(Y) - T(Y)| \geq c \} \leq \text{MSE}[t(Y)] / c^2$$

EXISTENCE RESULTS

$\int t^*$ IS UMVUE: $E_{p(s)} t^* = T(Y), \forall t \quad \forall_{p(s)} (t^*) \in \mathcal{V}_{p(s)} t$

$\int \Omega$ BE THE RANGE OF Y : $\Omega = \left\{ \prod_{i=1}^N (a_i, b_i) : a_i < y_i < b_i \right\}$

$\forall A \in \Omega$:

$$t_A = t^*[Y] - t^*[A] + A, \quad A = \sum_{i=1}^N A_i$$

$$E[t_A] = T(Y)$$

$$\forall [t_A]: E_{p(s)} [t^*[Y] - t^*[A] + A - Y]^2 = 0$$

WHEN $U = A$

$$\forall_p \{ t^* \} \subseteq \inf_{\Omega} \forall [t_A] = 0$$

\Rightarrow

ONLY DESIGN / ESTIMATOR WITH 0 VARIANCE
CAN BE UMVUE

GODAMBE (1955)

LINEAR ESTIMATORS:

$$t = \sum b_{is} y_i I_{js}$$

\Rightarrow

$t \equiv t_{HT}$ WITH OTHER DESIGN RESTRICTIONS

MOMENT GENERATING FUNCTIONS

ASSUME $T(Y) = 0$

$$M(t) = \mathbb{E} \exp \left[t \sum_{j \in S} Y_j \right] = \mathbb{E}_{P(S)} \exp \left[\sum_{j=1}^N t I_{jS} Y_j \right] =$$

$$= \{ \text{SRS} \} = \frac{1}{C_N^n} \sum_{S: n(S)=n} \exp \left[t \sum_{j \in S} Y_j \right]$$

CONDITIONAL OF BERNULLI SAMPLE:

$$P_{\uparrow \text{ IN THE BASK}} = P[n(S)=n] = C_N^n \lambda^n (1-\lambda)^{N-n}, \quad \lambda = n/N$$

JOINT MGF OF SAMPLE SIZE AND SAMPLE SUM

$$\begin{aligned} M(u, t) &= \mathbb{E} \exp \left[\sum_{j=1}^N t I_{jS} Y_j + I_{jS} u \right] = \prod_{j=1}^N \mathbb{E} \exp [I_{jS} (t Y_j + u)] \\ &= \prod_{j=1}^N [\lambda e^{t Y_j + u} + (1-\lambda)] \end{aligned}$$

EXPECTATIONS SUBTRACTED:

$$\left. \begin{aligned} n(S) - n &= \sum_{j=1}^N (I_{jS} - \lambda) \\ \sum_{j \in S} Y_j - \mathbb{E} \sum_{j \in S} Y_j &= \sum_{j=1}^N (I_{jS} - \lambda) Y_j \end{aligned} \right\}$$

$$\begin{aligned} M(u, t) &= \mathbb{E} \exp \left[\sum_{j=1}^N t (I_{jS} - \lambda) Y_j + (I_{jS} - \lambda) u \right] = \\ &= \prod_{j=1}^N [\lambda e^{(1-\lambda)(t Y_j + u)} + (1-\lambda) e^{-\lambda(t Y_j + u)}] \end{aligned}$$

COND'L MGF: DIVIDE BY PROB $n(S)=n$

$$M(t) = \prod_{j=1}^N [\lambda e^{(1-\lambda)t Y_j} + (1-\lambda) e^{-\lambda t Y_j}]$$

B

CENTRAL LIMIT THEOREM

SEQ OF ARRAYS

$$y_1^{(N)} \dots y_N^{(N)}, N \rightarrow \infty$$

$n =$ FUNCTION OF N

$$\frac{N-1}{N} \sqrt{p(s)} \sum_{j \in S} y_j = \underbrace{N \lambda (1-\lambda)}_n \underbrace{\sum_{j=1}^N (y_j - \mu_j)^2}_{\tilde{y}_j^2} / N =: V_N$$

TAILS CONDITIONS:

LINDBERGF: $\lim_{N \rightarrow \infty} d_N(\varepsilon) = 0 \quad \forall \varepsilon > 0$

$$d_N(\varepsilon) = \frac{1}{S_N^2} \sum_{j: |\tilde{y}_j| > \varepsilon \sqrt{V_N}} \tilde{y}_j^2$$

LYAPUNOV: $\lim_{N \rightarrow \infty} \frac{Q}{V_N^{2+\delta}} = 0$, $Q = \lambda(1-\lambda) \sum_{j=1}^N |y_j|^{2+\delta}$

$$V = N \lambda (1-\lambda) = n (1-n/N) \rightarrow \infty \text{ as } N \rightarrow \infty$$

CLT: APPROXIMATE $M(t)$ NEAR ZERO,

IGNORE CONTRIBUTIONS OF 'OUTLIER' TERMS

$$\Rightarrow \sum_{j \in S} -n \mu_j / V_N^{1/2} \xrightarrow{d} N(0,1)$$

MORE COMPLEX DESIGNS —

DEPENDS ON THE SAMPLING SCHEME:

- SIMPLE REJECTIVE DESIGN — HAJEK (1964)
- SUCCESSIVE SAMPLING — ROSEN (1972)
- STRATIFIED SAMPLES — BICKEL & FREEDMAN (1984)
($k \rightarrow \infty$ OR $N_h \rightarrow \infty$, $2 \leq n_h \leq N_h - 1$)
- STRATIFIED SAMPLES — KREWSKI & KAO (1981)
 n_h BDD, $k \rightarrow \infty$, $\max W_h = O(H^{-1})$
+ LYAPUNOV CONDⁿ, + CONDⁿ OF POPⁿ COVⁿ MX
- HORVITZ - THOMPSON : P.K. SEN (1980, 1988)
- MULTIVARIATE EXTENSIONS

DESIGN-CONSISTENCY

WE KNOW: $\frac{e - T(y)}{\sqrt{V(e)}} \xrightarrow{d} N(0,1)$

$T(y)$ AND $V(e)$ (ESP THE LATTER) ARE DESIGN-BASED CONCEPTS.

DOES THE STUDENTIZED VERSION $\frac{e - T(y)}{\sqrt{v(e)}}$ ALSO $\xrightarrow{d} N(0,1)$?

$v(e)$ HAS TO BE DESIGN-CONSISTENT!

$$\frac{v(e)}{V(e)} \xrightarrow{p} 1$$

HERE, p IS THE PROBⁿ OVER SAMPLES $p(s)$
HENCE, NEED TO BE ABLE TO SHOW
CONSISTENCY OF $v(e)$

FOR THE BASIC DESIGNS CONSIDERED THUS
FAR (SRS, STRATIFIED, HT + YGS, MULTI-STAGE)
AND BASIC $v(e)$, THIS CONSISTENCY DOES HOLD

SVY - AUX VARIABLES

Note Title

2/6/2007

SO FAR: FOCUS ON SINGLE VARIABLE ESTIMATORS
SIMPLE EXTENSIONS: AUX INFO ESTIMATORS

RATIO ESTIMATOR

SUPPOSE VARS Y AND X ARE COLLECTED

$$(y_j, x_j), j \in S$$

INTEREST IS IN RATIO $R = T[Y]/T[X]$

\Rightarrow ESTIMATE BY

$$r = t[y]/t[x] \quad (\text{COMBINED RATIO ESTIMATOR})$$

OR

$$r_s = \sum_{h=1}^H W_h r_h, \quad r_h = t_h[y]/t_h[x] \\ (\text{SEPARATE RATIO ESTIMATOR})$$

RATIO ESTIMATOR OF TOTAL:

$$t_r[y] = r T[x] \\ \uparrow \text{ASSUMED KNOWN}$$

REGRESSION ESTIMATOR OF TOTAL:

$$t_L[y] = t[y] + b(T[x] - t[x]), \\ b = \frac{S_{xy}}{S_x^2} = \text{SAMPLE REGRESSION COEFF}$$

$$S_{z\eta} = \frac{1}{n-1} \sum_{j=1}^n (z_j - \bar{z})(\eta_j - \bar{\eta}), \quad z, \eta = X, Y$$

$$S_{z\eta} = \frac{1}{n-1} \sum_{j \in S} (z_j - \bar{z})(\eta_j - \bar{\eta}), \quad z, \eta = X, Y$$

SAMPLE AVERAGE

(OR DESIGN-BASED ESTIMATOR, IN GENERAL)

COEFF OF VARIATION: $C_{z\eta} = S_{z\eta} / \bar{z}\bar{\eta}$

$E r \neq R \Leftrightarrow$ NON-LINEARITY

BIAS $B(r) = ?$

VARIANCE $V(r) = ?$

IF THE LATTER IS UNAVAILABLE, MSE $M(r) = ?$

LINEARIZATION / TAYLOR SERIES EXPANSION APPROACH

SUPPOSE $\theta = f(T(Y))$ IS A POP ~ PARAM

$\hat{\theta} = f(t(Y))$ IS A CONSISTENT ESTIMATOR

$$\hat{\theta} - \theta = \nabla f \cdot (t(Y) - T(Y)) + \frac{1}{2} (t(Y) - T(Y))' [\nabla^2 f] (t(Y) - T(Y)) + \text{SMALLER ORDER TERMS}$$

SUPPOSE $E_{p(Y)} t(Y) = T(Y) \Rightarrow$

$$B \hat{\theta} = \frac{1}{2} E (t(Y) - T(Y))' [\nabla^2 f] (t(Y) - T(Y)) \approx O_p(n^{-1})$$

$$M(\hat{\theta}) = E(\hat{\theta} - \theta)^2 \approx \nabla f' V(t(Y)) \nabla f \quad (*)$$

ON THE OTHER HAND,

$$\begin{aligned}\hat{\theta} - \theta &\approx \nabla f' \left(\sum_{j \in S} \frac{y_j}{n_j} - T[Y] \right) \\ &= \sum_{j \in S} \left\{ \nabla f' y_j \right\} / n_j - \nabla f' T[Y] \\ &= \text{A LINEAR STATISTIC OF} \\ &u_j = \nabla f' y_j\end{aligned}$$

$$\text{So } M[\hat{\theta}] \approx \sum_{j \neq k} \left(\frac{u_j}{n_j} - \frac{u_k}{n_k} \right)^2 \rho_{jk}$$

$$m[\hat{\theta}] \approx \sum_{\substack{j \neq k \\ \in S}} \left(\frac{u_j}{n_j} - \frac{u_k}{n_k} \right)^2 \frac{n_j n_k - n_{jk}}{n_{jk}}$$

BACK TO THE RATIO ESTIMATOR

DESIGN = SRS

$$B[t_r(Y)] = \frac{1-f}{n} \bar{y}^2 C_x(C_x - \rho(Y)) \quad , \quad (C_x = \sqrt{C_{33}})$$

$$M[t_r(Y)] \approx V[t_r(Y)] = \frac{1-f}{n} S_d^2$$

$$S_d^2 = \sum_{j=1}^N \frac{(y_j - R x_j)^2}{N-1} = S_{xy}^2 - 2RS_{xy} + R^2 S_{xx} \quad (**)$$

WITH PLUG-IN ESTIMATOR

$M[t_r] \leq V[t_r]$ FOR REASONABLE DESIGNS

$$\left| B(t_r(Y)) / V^{1/2}(t_r(Y)) \right| \leq C_{\bar{x}} = \sqrt{(1-f) C_{xx} / n}$$

APPLICATION: SAMPLE MEAN IN CLUSTERED SAMPLES

DESIGN: 1ST STAGE = SRS OF PSUs, l OUT OF L

2ND STAGE = SRS OF $n_r = \sum N_r$ UNITS

$$\bar{y} = \frac{t[y]}{t[1]} = \frac{\sum_{r=1}^l N_r \bar{y}_r}{\sum N_r} = \frac{\sum_r \sum_{j \in B_r} y_j}{\sum_r n_r} = \frac{\sum y}{n}$$

$$\sigma_r = \sum_{j \in B_r} y_j$$

$$s(\bar{\sigma}) = \frac{1}{l(l-1)} \sum_r (\sigma_r - \bar{\sigma})^2$$

$$s(\bar{n}) = \frac{1}{l(l-1)} \sum_r (n_r - \bar{n})^2$$

$$\text{cov}(\bar{\sigma}, \bar{n}) = \frac{1}{l(l-1)} \sum_r (\sigma_r - \bar{\sigma})(n_r - \bar{n})$$

$$\text{var}(\bar{y}) = \left(\frac{1}{\bar{n}^2} \right) \left[s(\bar{\sigma}) - 2\bar{y} \text{cov}(\bar{\sigma}, \bar{n}) + \bar{y}^2 s(\bar{n}) \right] =$$

$$= \frac{1}{\bar{n}^2} \frac{1}{l(l-1)} \sum_r \left[(\sigma_r - \bar{\sigma})^2 - 2\frac{\bar{\sigma}}{\bar{n}} (\sigma_r - \bar{\sigma})(n_r - \bar{n}) + \frac{\bar{\sigma}^2}{\bar{n}^2} (n_r - \bar{n})^2 \right] =$$

$$= \frac{1}{\bar{n}^2 l(l-1)} \sum_r \left[(\sigma_r - \bar{\sigma}) - \frac{\bar{\sigma}}{\bar{n}} (n_r - \bar{n}) \right]^2 =$$

$$= \frac{1}{\bar{n}^2 l(l-1)} \sum_r n_r \left[\frac{\sigma_r}{n_r} - \bar{y} \right]^2$$

" σ_r/n_r "

A STRANGE
LOOKING
ESTIMATOR!

HOMEWORK:

- (i) DERIVE THE VARIANCE OF THE RATIO (x/y) USING (*)
- (ii) FIND THE LINEARIZED CONTRIBUTIONS U_r FOR THE RATIO ESTIMATOR, AND FIND THEIR VARIANCE AND ITS ESTIMATOR

JG 2.4.9

REGRESSION ESTIMATOR

$$t_{\text{reg}}[y] = t[y] + b(\bar{T}[x] - t[x])$$
$$V[t_{\text{reg}}] \approx M[t_{\text{reg}}] \approx \frac{1-f}{n} \frac{1}{N-1} \sum_{j=1}^N \Delta_j^2 = \frac{1-f}{n} S_y^2 (1-\rho^2)$$

\uparrow
SRS

$$\Delta_j = (Y_j - \bar{Y}) - B(X_j - \bar{X}) \quad \leftarrow \text{LINEAR TERM}$$

$B =$ POPULATION REGRESSION COEFF

MOST OF THE TIME,

$$V[t_{\text{reg}}] \leq V[t_{\text{rat}}] \leq V[t_{HT}]$$

\uparrow

FOR REASONABLY STRONGLY
CORRELATED X, Y

OTHER ESTIMATORS OF VARIANCE:
MODEL-BASED / ASSISTED

SVY - REGRESSION W/ COMPLEX DATA

Note Title

2/11/2007

RESOURCES:

BINDER (1983)

SKINNER (1989)

FULLER (2002)

I.I.D.:

$$Y_j = X_j \beta + e_j$$

$$E e_j = 0$$

$$\text{OLS} : \sum_{j=1}^N (y_j - x_j \beta)^2 \rightarrow \min$$

- IS THIS A SENSIBLE THING?

MAY NOT BE A CORRECT MODEL,
BUT AT LEAST THAT'S SOMETHING
YOU CAN ALWAYS COMPUTE

NORMAL ERRORS:

$$\sum_{j=1}^N x_j (y_j - x_j \beta) = 0$$

- A POPULATION TOTAL OF ZERO $E R^1$!

SOMETHING WE COULD TRY TO
ESTIMATE

SAMPLE / ESTIMATOR:

$$\sum_{j \in S} x_j (y_j - x_j b) = 0 \quad \leftarrow \text{OLS}$$

$$\sum_{j \in S} \frac{x_j (y_j - x_j b)}{\pi_j} = 0 \quad \leftarrow \text{HT}$$

ANALOGY
TO
RATIO
ESTIM.

$$b_w = (X' W X)^{-1} (X' W Y)$$

$$W = \text{DIAG}(\pi_j^{-1})$$

UNLESS THE DESIGN IS SELF-WEIGHTING,
THE OLS ESTIMATOR MIGHT BE
BIASED.

HOWEVER, IF

...

BLAH-BLAH ABOUT $\pi_j^{-1} \in \mathcal{L}(X)$

SEE
LATER

THUS TO ENSURE DESIGN CONSISTENCY,

- ① USE WEIGHTS
- ② USE DESIGN INFO IN THE MODEL

LET US REVISIT OLS $b_1 = b_2$ AGAIN

$$b_1 = \left(\sum_{j \in S} x_j x_j' \right)^{-1} \left(\sum_{j \in S} x_j' y_j \right)$$

SUPPOSE THE DESIGN IS SUCH THAT $b_1 \xrightarrow{P} B$

$$b_1 - B = ?$$

RECALL THE LINEAR FRAMEWORK:

- 1) BASIC ESTIMATOR t
- 2) TAYLOR SERIES LINEAR TERM ∇f
- 3) GET $V[t]$
- 4) $M(\hat{\theta}) \approx \nabla f \cdot V[t] \cdot \nabla f'$

1) BASIC ESTIMATOR:

$$t_b = \sum_{j \in S} x_j (y_j - x_j' b) = 0$$

$$T = \sum_{j=1}^N x_j (y_j - x_j' B) = 0$$

$$2) \nabla f(b) = \frac{\partial T}{\partial B} = \sum_{j \in S} x_j x_j', \quad E \nabla f(b) = \frac{n}{N} \sum_{j=1}^N x_j x_j'$$

$$\begin{aligned} 3) V[t] &= E_{p(S)} t_b t_b' = E \left(\sum_{j \in S} x_j e_j \right) \left(\sum_{j \in S} x_j e_j \right)' \\ &= E \sum_{j \in S} x_j x_j' e_j^2 + E \sum_{j \neq k} x_j x_k' e_j e_k \end{aligned}$$

$$\text{SRS: } \frac{n}{N} \sum_j^N x_j x_j' e_j^2 - \frac{n(n-1)}{N(N-1)} \sum_{j=1}^N x_j x_j' e_j^2$$

$$\sum_{k \neq j} x_k e_k = -x_j e_j$$

$$= \frac{n(N-1) - n(n-1)}{N(N-1)} \sum_j x_j x_j' e_j^2$$

$$= \frac{n(N-n)}{N(N-1)} \sum_j x_j x_j' e_j^2$$

$$\begin{aligned} v(b) &\approx \left(\frac{n}{N} \sum_{j=1}^N x_j x_j' \right)^{-1} \frac{n(N-n)}{N(N-1)} \sum_j x_j x_j' e_j^2 \left(\frac{n}{N} \sum_{j=1}^N x_j x_j' \right)^{-1} \\ &\approx \frac{1}{n} \left(\frac{1}{N} \sum_{j=1}^N x_j x_j' \right)^{-1} \left[\frac{1}{N} \sum_j x_j x_j' e_j^2 \right] \left(\frac{1}{N} \sum_{j=1}^N x_j x_j' \right)^{-1} \end{aligned}$$

$$v(b) \approx \frac{1}{n} \left(\frac{1}{n} \sum_{j \in S} x_j x_j' \right)^{-1} \left(\frac{1}{n} \sum_{j \in S} x_j x_j' e_j^2 \right) \left(\frac{1}{n} \sum_{j \in S} x_j x_j' \right)^{-1}$$

AKA SANDWICH/ROBUST ESTIMATOR

$$\text{COMPLEX DESIGN: } v_L(b) = v_{xx}^{-1} v_{xy} v_{xx}^{-1}$$

$$v_{xx} = \frac{1}{n} \sum_h \sum_c \sum_j x_{hej} x_{hej}^T$$

$$v_{xy} = \hat{V} \left[(y - X\hat{b})x \right] = \sum_h w_h \frac{(1-f_h)e_h}{l_h-1} \sum_e (u_{he} - \bar{u}_h) (u_{he} - \bar{u}_h)^T$$

$$u_{he} = \sum_{j \in \mathcal{B}_e \cap \mathcal{S}_h} u_{hej}, \quad u_{hej} = \frac{1}{n} \cdot x_{hej} \left[y_{hej} - \hat{b}^T x_{hej} \right]$$

USE OF DESIGN INFO

$$\bar{y}_{\text{reg}} = \bar{X}_N \hat{\beta}$$

KNOWN POP ~ TOTAL

CALIBR ~:

$$\sum_{j \in S} d_j \alpha_j = T[\alpha]$$

$$E_{p(s)} \bar{X}_N \hat{\beta} = E_{p(s)} \bar{X}_N (X'WX)^{-1} X'WY$$

$$(X'WX)^{-1} X'WY = (X'WX)^{-1} (X'WX\beta + X'WE) = \beta + (X'WX)^{-1} X'WE$$

FOR IT TO BE UNBIASED, NEED TO HAVE

$$E_{p(s)} (X'WX)^{-1} X'WE = 0$$

SUFFICIENT COND ~: $E_{p(s)} X'WE = 0$

LOCATION INVARIANCE WRT Y : $W D_{\pi}^{-1} \mathbb{1} \in \mathcal{L}(X)$

SO THAT THE DESIGN EXPECTN
CAN BE CARRIED THROUGH $X'W$
(OTHERWISE, Y 'S AND E 'S MAY "GET LOST"
IN THE X^{\perp} SPACE, AND $E_{p(s)}$ WILL
PRODUCE BIASES)

WAYS TO ENSURE THAT:

1. HAVE $\mathbb{1} \in X$, USE $W = D_{\pi}$

\Rightarrow SUPPOSE THIS IS THE SINGLE REGRESSOR, THEN

$$\bar{y} = \left(\sum_{j \in S} \pi_j^{-1} \right)^{-1} \sum_{j \in S} \pi_j^{-1} y_j = t_{HT}[Y] / t_{HT}[\mathbb{1}]$$

2. HAVE $W = I$, AND USE π_j^{-1} AS ONE OF REGRESSORS

3. HAVE π_j AS ONE OF REGRESSORS, AND
USE $W = D_{\pi}^{-2}$
 $\Rightarrow \bar{y} = N^{-1} \sum_{j \in S} y_j / \pi_j$ - HT ESTIMATOR

4. STRATA WITH FRACTIONS f_h , SRS WITHIN STRATA \Rightarrow CAN USE STRATA INDICATORS TO ENSURE THE LINEAR SPACE RESTRICTION

\hookrightarrow AN EXAMPLE OF MODELLING DESIGN RATHER THAN PROVIDING DESIGN-BASED INFERENCE

MULTIVARIATE DESIGN EFFECTS (SEC. 2.11 OF SRS)

VECTOR $\hat{\beta} \in \mathbb{R}^p \Rightarrow$ WHAT IS A PROPER WAY TO LOOK AT DEFF?

1D: $DEFF(t, v_0) = \frac{V_{p(s)}[t]}{V_0[t]} \leftarrow$ ESTIMATOR UNDER SOME "NULL" ASSUMPTIONS

$$DEFF(\hat{\theta}, v_0) = \Delta = E(v_0^{-1}) V(\hat{\theta})$$

TYPICALLY, $v_0 = V_{SRS}[\hat{\theta}]$

GENERALIZED DESIGN EFFECTS:

$$\max_{c \neq 0} \text{def}(c'\hat{\theta}, c'V_0c) = \delta_1 \geq \dots \geq \delta_p = \min_{c \neq 0} (c'\hat{\theta}, c'V_0c) = \text{Spec}(\Delta)$$

EFFECT OF GENERALIZED DEFFS:

$$\chi^2 = (\hat{\theta} - \theta)' V_0^{-1} (\hat{\theta} - \theta) \sim \sum_{i=1}^p \delta_i \chi_{1,i}^2$$

(RAO & SCOTT, 1981)

SATTERTHWAITE (1946) APPROX~:

$$\frac{\chi^2}{\delta(1+a^2)} \sim \chi_q^2 \quad q = \frac{p}{1+a^2}, \quad a = \text{CV}\{\delta\}$$

↑
COEFF OF VARIATION

$$p\bar{\delta} = \text{tr} \Delta, \quad p\bar{\delta}^2(1+a^2) = \text{tr} \Delta^2$$

MISSPECIFICATION EFFECT

(SEC 2.2 OF SHS)

- EFFECT OF THE DESIGN ON VARIANCE ESTIMATORS

$$\text{MEFF}(\hat{\theta}, V_0) = \frac{V_{\text{true}}[\hat{\theta}]}{E_{\text{true}}[V_0]}$$

MAY BE DESIGN-BASED OR MODEL-BASED

CLUSTER DEFF ($n=2$) : $1 + \rho_{\text{ICC}}$ - VARIANCE IS ρ_{ICC} % MORE THAN WHAT IT WOULD BE UNDER SRS

CLUSTER MEFF : $\frac{1 + \rho_{\text{ICC}}}{1 - \rho_{\text{ICC}}}$ - VARIANCE IS 3 TIMES MORE THAN WHAT AN SRS-BASED ESTIMATOR WILL TELL

OFTEN: $E_{\text{TRUE}} \sigma_0 \approx V_{\text{SRS}}[\hat{\theta}]$

BOTH ARE USUALLY ESTIMATED BY SAMPLE-BASED σ

DESIGN EFFECTS IN REGRESSION

(SHS, SEC. 3.3.4)

WHAT IS THE APPROPRIATE σ_0 ?

(i) SRS + HOMOSKEDASTICITY: $\sigma_{OLS} = s_e^2 (X'X)^{-1}$

(ii) SRS + HETEROSKEDASTICITY: SANDWICH

(SEE ALSO OTHER CORRECTIONS IN MY
MISSPECIFICATION TALK)

MEFF $(\hat{\beta}_{OLS}, V_{OLS}) \approx 1 + \rho C_\sigma C_x$

$C_\sigma = CV \sigma^2(x_j)$, $C_x = CV x_j$, $\rho = \text{corr}(\sigma^2(x_j), x_j)$

(iii) NESTED ERRORS / GLS:

$$Y_{hj} = x'_{hj} \beta + u_h + e_{hj}$$

$$E e_{hj} e_{h'j'} = \sigma_u^2 \delta_{hh'} + \sigma_e^2 \delta_{hh'} \delta_{jj'}$$

$$\frac{V_{\text{NESTED ERRORS}}(\text{SLOPE})}{V_{\text{OLS}}(\text{SLOPE})} \approx 1 + \underset{\substack{\uparrow \\ \text{CLUSTER} \\ \text{SIZE}}}{(M-1)} \rho_{\text{corr}, \varepsilon} \rho_{\text{corr}, x}$$

TYPICAL DEFFs:

MEAN $\rightarrow 2$

SLOPE $\rightarrow 1.2$

CORR $\sim \rightarrow 2$

$R^2 \rightarrow 5$

RECALL THAT REGRESSING ON DESIGN INFO
HELPS CONSISTENCY... REGRESSING ON VARS
CORRELATED WITH DESIGN INFO (E.G. x 'S THAT
ARE CORRELATED WITHIN CLUSTERS IN THE SAME
WAY AS y 'S) ACTS AS A SURROGATE

SVY - NONLINEAR MODELS

Note Title

2/17/2007

BINDER (1983)
SKINNER (1989)
THOMPSON (1997, ch. 4, 6)

SEC. 4.1

POPULATION ESTIM~ EQ:

$$\sum_{j=1}^N \varphi_j(y_j, x_j, \theta_N) = 0$$

MEAN: $\varphi_j = y_j - \theta_N$
RATIO: $\varphi_j = y_j - \theta_N x_j$
C.D.F.: $\varphi_j = \mathbb{1}\{y_j \leq y\} - \theta_N$

$$\varphi_s(\theta) = \sum_{j \in S} \varphi_j(y_j, x_j, \theta) / \pi_j$$

$$E_{p(s)} \varphi_s(\theta) = \sum_{j=1}^N \varphi_j(y_j, x_j, \theta)$$

$$V_{p(s)} \varphi_s(\theta) = \sum_{j=1}^N \left(\frac{\varphi_j}{\pi_j} \right)^2 \pi_j (1 - \pi_j) + \sum_{j \neq k} \frac{\varphi_j \varphi_k}{\pi_j \pi_k} (\pi_{jk} - \pi_j \pi_k)$$

$$\text{ESTIMATOR: } \sigma(\varphi_s) = \frac{1}{2} \sum_{j \neq k} w_{jk} \left(\frac{\varphi_j}{\pi_j} - \frac{\varphi_k}{\pi_k} \right)^2$$

$$\text{CLT: } \left(\varphi_s(\theta) - \sum_{j=1}^N \varphi_j(y_j, x_j, \theta) \right) / \sqrt{V_{p(s)}[\varphi_s(\theta)]} \xrightarrow{d} N(0, 1)$$

$w_{jk} = (\pi_j \pi_k - \pi_{jk}) / \pi_{jk}$

WITH APPROPRIATE STANDARDIZED VERSION

TAYLOR SERIES:

$$\frac{1}{N} \varphi_s(\hat{\theta}_s) - \frac{1}{N} \varphi_s(\theta) = -\frac{1}{N} \varphi_s(\theta) =$$

$$= \frac{\partial \varphi}{\partial \theta} (\hat{\theta}_s - \theta) + \frac{1}{2} (\hat{\theta}_s - \theta)' \frac{\partial^2 \varphi}{\partial \theta \partial \theta} \Big|_{\bar{\theta}} (\hat{\theta}_s - \theta)$$

$\bar{\theta}$ - INTERMEDIATE POINT

SEC. 4.2 - FUNCTIONS OF TOTALS, LINEARIZATION TECHNIQUE

$$\theta = g(T), \quad \hat{\theta} = g(t)$$

$$\begin{aligned} \Rightarrow \varepsilon_L(\hat{\theta}; \theta) &= \hat{\theta} - \theta = g(t) - g(T) = \\ &= \sum_{\alpha} \frac{\partial g}{\partial T_{\alpha}} \Big|_{T_{\alpha}} (t_{\alpha} - T_{\alpha}) + o(\dots) \end{aligned}$$

@ POP ~ TOTALS

$$\begin{aligned} n^{1/2}(t - T) &\sim N(0, \dots), \quad \left\| \frac{\partial g}{\partial T} \right\| > 0 \\ \Rightarrow \frac{g(t) - g(T)}{V^{1/2} \left\{ \sum_{\alpha} \frac{\partial g}{\partial T_{\alpha}} \Big|_{T_{\alpha}} \right\}} &= \frac{\hat{\theta} - \theta}{V^{1/2}[\varepsilon]} \xrightarrow{d} N(0, 1) \end{aligned}$$

ESTIMATOR: $\sigma(\varepsilon_L) = \sigma \left(\sum_{j \in S} u_j / \pi_j \right) \Big|_t$

E.G. Y-G-S ESTIMATOR

$$u_j = \sum_{\alpha} \frac{\partial g}{\partial T_{\alpha}} Y_{\alpha j}$$

4.2.3 - 4.2.6:

RANDOM GROUPS, BRR, JACKKNIFE, BOOTSTRAP

SEC. 4.3:

U-STATISTICS WITH SRS

SEC. 4.4:

ESTIMATING EQNS WITH
NUISANCE PARAMETERS

MOST IMPORTANT EXAMPLE -

THE COX PROP HAZARD MODEL

BINDER (1983)

DESIGN-BASED SAMPLING DESIGN
FOR PARAMETERS DEFINED AS
FUNCTIONS OF DATA VALUES

E.G. REGRESSION: $X^T X B = X^T Y$
GENERALIZED LINEAR MODEL

↑ REVIEW OF GLM:

$$Y \sim f(y; \theta, \varphi) = \exp \left[\underset{\uparrow}{\alpha(\varphi)} \{y\theta - g(\theta) + h(\varphi)\} + \delta(\varphi, y) \right]$$

MORE OFTEN GOES INTO DENOMINATOR

$$E Y = g'(\theta) \equiv \mu(\theta) \quad - \text{MEAN FN}$$

$$V Y = \mu'(\theta) \alpha(\varphi) \quad - \text{VARIANCE FN}$$

$$\theta = f(x; \beta) \quad - f \text{ IS A KNOWN FN}$$

β IS A VECTOR OF COEFFTS

SCORE EQ:

$$\sum [Y_k - \mu_k(f(x; \beta))] f'(x; \beta) x = 0$$

REGRESSION: $f = \text{IDENTITY LINK}$, $g = \mu^2/2$, $\alpha = \sigma^2$

LOGIT: $f = \text{IDENTITY}$, $g = \ln p/(1-p)$, $\alpha = 1$
etc.

↓ IMPLICIT PARAMETERS:

$$\text{REGRESSION: } X^T X B = X^T Y$$

$$\text{GLM: } \sum_{j=1}^N [Y_j - \mu\{f(x_j^T B)\}] f'(x_j^T B) x_j = 0$$

$$\text{GENERALLY: } \sum_{j=1}^N \varphi_j(Y_j; x_j; B) = 0$$

$$\text{ESTIM~ FRAMEWORK: } W_N(\theta) = \sum_{j=1}^N \varphi_j(z_j; \theta) - V(\theta)$$

POP~ VALUE: $W_N(\theta_0) = 0$

REGULARITY COND~:

- (i) $\theta \in \text{int } \Theta$; NBHOD $U(\theta) \in \Theta$
- (ii) SEQ OF POP~ AND DESIGNS: $\{U_N, P_N(s)\}$:
 - SCALED $(t-T) \xrightarrow{d} N(0, \Sigma)$
 - $U(t)/V(t) \xrightarrow{p} 1$
- (iii) CONTINUITY & LIMITS OF $w_N(\theta), \partial w_N(\theta)$
- (iv) CONTINUITY OF $V(t)$

ESTIM~ OF θ :

ASSUME $t[U(\theta)]$ IS ASYMPT~ NORMAL:
 SCALING $(t[U(\theta)] - U(\theta)) \xrightarrow{d} N(0, \Sigma_u(\theta))$
 ESTIMATED BY $\hat{\Sigma}_u(\theta)$

$\hat{w}(\theta) = t[U(\theta)] - v(\theta)$ - ESTIM EQU FOR θ :
 $\hat{\theta}$ IS SUCH THAT $\hat{w}(\hat{\theta}) = 0$

REGRESSION: $u(y_j; x_j; \beta) = -(y_j - x_j^T \beta) x_j, v(\beta) = 0$

$$0 = \hat{w}(\hat{\theta}) \approx \hat{w}(\theta_0) + \left. \frac{\partial \hat{w}(\theta)}{\partial \theta} \right|_{\theta_0} (\hat{\theta} - \theta_0) + \text{REMAINDER}$$

$$\hat{w}(\theta_0) \approx - \frac{\partial \hat{w}(\theta_0)}{\partial \theta} (\hat{\theta} - \theta_0)$$

$V[\hat{\theta}] \rightarrow$

$$\Sigma_u(\theta_0) \approx \frac{\partial w_N(\theta_0)}{\partial \theta} V(\hat{\theta}) \left[\frac{\partial w_N(\theta_0)}{\partial \theta} \right]^T$$

$$\Rightarrow V[\hat{\theta}] = \left[\frac{\partial w_N(\theta_0)}{\partial \theta} \right]^{-1} \Sigma_u(\theta_0) \left[\frac{\partial w_N(\theta_0)}{\partial \theta} \right]^{-T}$$

WITH ESTIMATOR

$$V[\hat{\theta}] = \left[\frac{\partial w(\hat{\theta})}{\partial \theta} \right]^{-1} \hat{\Sigma}_u(\hat{\theta}) \left[\frac{\partial w(\hat{\theta})}{\partial \theta} \right]^{-T}$$

GLM, CANONICAL LINK $f(\cdot) = \text{IDENTITY}$:

$$\Rightarrow \sum_{j=1}^N [Y_j - \mu(x_j^T \beta)] X_j = 0$$

$$w_N(\beta) = \sum_{j=1}^N [Y_j - \mu(x_j^T \beta)] X_j$$

$$\frac{\partial w_N}{\partial \beta} = - \sum_{j=1}^N \mu'(x_j^T \beta) X_j X_j^T = X^T \Lambda(\beta) X$$

$$\Lambda(\beta) = \text{diag} [\mu'(x_1^T \beta), \dots, \mu'(x_N^T \beta)]$$

$$\Rightarrow V(\hat{\beta}) = [X^T \Lambda(\beta) X]^{-1} \Sigma(\beta) [X^T \Lambda(\beta) X]^{-1}$$

$$\Sigma(\beta) = V \left[t(e_j; x_j) \right]$$
$$e_j = Y_j - \mu(x_j^T \beta)$$

EACH MATRIX, AND EACH MATRIX ENTRY, IS EITHER A TOTAL OR AN ESTIMATE OF VARIANCE OF A TOTAL

$$\Rightarrow V(\hat{\beta}) = \hat{J}^{-1}(\hat{\beta}) \hat{\Sigma}(\hat{\beta}) \hat{J}^{-1}(\hat{\beta})$$

WHERE

$$\hat{J}(\hat{\beta}) = t [X_j X_j^T \mu'(x_j^T \hat{\beta})]$$

(AKA DERIVATIVES / JACOBIAN

IN NEWTON-RAPHSON STEPS)

REGRESSION: $w_N(\beta) = X^T Y - X^T X \beta$; $\frac{\partial w_N(\beta)}{\partial \beta} = -X^T X$

$$\Rightarrow V(\hat{\beta}) = S_{XX}^{-1} \hat{\Sigma}(\hat{\beta}) S_{XX}^{-1}$$

$$S_{XX} = t(X_j X_j^T) \quad \hat{\Sigma}(\hat{\beta}) = V(\hat{e}_j; x_j)$$

EXTENSION:

ESTIMATE SIMULTANEOUSLY $(\hat{\gamma}, \hat{B}, R^2)$:

$$w_N(\theta) = \begin{cases} Y^T \mathbb{1} - N\bar{Y} \\ X^T Y - X^T X B \\ (Y^T Y - N\bar{Y}^2)(R^2 - 1) + Y^T Y - Y^T X B \end{cases} = \begin{cases} 0 \\ 0 \\ 0 \end{cases}$$

$$\hat{R}^2 = 1 - \frac{S_{yy} - \hat{B}^T S_{xy}}{S_{yy} - N\bar{y}^2}$$

$$\frac{\partial w}{\partial \theta} = \begin{pmatrix} -N & 0 & 0 \\ 0 & -X^T X & 0 \\ 2N\bar{y}(1-R^2) & -Y^T X & Y^T Y - N\bar{y}^2 \end{pmatrix} \rightarrow \left(\frac{\partial w}{\partial \theta} \right)^{-1}$$

$$u_j = \begin{pmatrix} y_j \\ (y_j - x_j^T \hat{B}) x_j \\ (R^2 y_j - x_j^T \hat{B}) y_j \end{pmatrix} \rightarrow \sigma[u_j]$$

$$\Rightarrow \text{Var}[R^2] = (S_{yy} - N\bar{y}^2) a^T \text{Var}[u] a$$

$$a = 2\bar{y}(1-R^2) - \hat{B}^T \mathbb{1}$$

LOGISTIC REGRESSION:

$$\mu(x^T B) = \frac{\exp(x^T B)}{1 + \exp(x^T B)}$$

$$w_N(B) = \sum_{j=1}^N \left[y_j - \frac{\exp(x_j^T B)}{1 + \exp(x_j^T B)} \right] x_j$$

$$X^T N(B) X = \sum_{j=1}^N x_j x_j^T \frac{\exp(x_j^T B)}{[1 + \exp(x_j^T B)]^2}$$

$$\hat{\Sigma}(\hat{B}) = \sum_{j \in S} x_j x_j^T \left[y_j - \frac{\exp(x_j^T B)}{1 + \exp(x_j^T B)} \right]^2$$

ANOTHER EXAMPLE (SEC. 4.4):

LOG-LINEAR MODELS FOR CATEG. DATA
(MULTINOMIAL LOGISTIC MODEL)

MOVE ON TO SKINNER (1989)

= CH 3 OF SHS (1989)

$$l(\theta) = \sum_j \ln f(y_j; \theta)$$

$$\hookrightarrow u_j(\theta) = \nabla \ln f(y_j; \theta)$$

$$\text{POP} \sim \text{PARR: } \theta_0 : T[u_j(\theta)] = 0$$

$$\text{PSEUDOMLE } \hat{\theta}_{\text{PML}} \leftarrow t[\hat{\theta}_{\text{PML}}] = 0$$

- NOT THE MLE ESTIMATOR
⇒ OPTIMALITY DOES NOT HOLD!
- VARIOUS ESTIMATORS t
⇒ DIFFERENT $\hat{\theta}_{\text{PML}}$!

$$\begin{aligned} \text{As. } \nabla [\hat{\theta}_{\text{PML}}] &= \text{LINEAR} \\ &= I(\hat{\theta}_{\text{PML}})^{-1} \nabla_L [t(\hat{\theta}_{\text{PML}})] I(\hat{\theta}_{\text{PML}})^{-1} \end{aligned}$$

$$I(\theta) = \frac{\partial}{\partial \theta} t[u(\theta)]$$

$\nabla_L [t(\hat{\theta}_{\text{PML}})] \leftarrow$ ESTIMATED IN ANY
REASONABLE DESIGN-CONSISTENT WAY

$$\begin{aligned} \text{LOGISTIC: } u_j(\theta) &= (y_j - p_j(\theta)) x_j \\ I(\theta) &= t[p_j(1-p_j) x_j x_j^T] \end{aligned}$$

WALD TESTS: χ^2 ; LR TESTS: GEN & DEFF

SVY: RESAMPLING ESTIMATORS

Note Title

3/14/2007

RESOURCES:

SHAO (1996) REVIEW

RAO & WU (1988) SCALED BOOTSTRAP

RAO, WU & YUE (1992) BOOTSTRAP IN WEIGHTS

KREWSKI & RAO (1981) CONSISTENCY RESULTS

RAO & WU (1985) COMPARISON OF ESTIMATORS

SUMMARY REVIEW: MY TALK, SEE:

<http://web.missouri.edu/~kolenikovs/talks/survey-resampling-by2.pdf>

BRR:

LET'S LOOK @ ESTIMATOR OF TOTAL.

$$\bar{y} = \sum_h W_h \bar{y}_h = \sum_h W_h \frac{y_{h1} + y_{h2}}{2}$$

$$\sigma(\bar{y}) = \sum_h \frac{W_h^2}{2} \sigma(\bar{y}_h) = \sum_h W_h^2 \left[\frac{(y_{h1} - \bar{y}_h)^2 + (y_{h2} - \bar{y}_h)^2}{2} \right]$$

$$= \sum_h \frac{W_h^2}{4} (y_{h1} - y_{h2})^2$$

DIVIDE THE SAMPLE INTO HALF-SAMPLES:

$$\bar{y}^{(1)} = \sum_h W_h y_{h1}, \quad \bar{y}^{(2)} = \sum_h W_h y_{h2}$$

$$\bar{y} = \frac{1}{2} (\bar{y}^{(1)} + \bar{y}^{(2)})$$

(NEARLY) INDEP \Rightarrow

$$\sigma(\bar{y}) = \frac{1}{4} (\bar{y}^{(1)} - \bar{y}^{(2)})^2$$

MORE EFFICIENT:

- TAKE ALL POSSIBLE $\frac{1}{2}$ -SAMPLES $\rightarrow 2^H$

$$\Rightarrow \frac{1}{2^H} \sum_{\nu} (\bar{y}^{(\nu)} - \bar{y})^2 = \sigma(\bar{y})$$

$$\begin{aligned} \delta_h^{(\nu)} &= \pm \frac{1}{2} \Delta y_h = \pm \frac{1}{2} (y_{h1} - y_{h2}) \\ \bar{y}^{(\nu)} - \bar{y} &= \bar{y} + \sum_h w_h \delta_h^{(\nu)} - \bar{y} = \sum_h w_h \delta_h^{(\nu)} \\ (\bar{y}^{(\nu)} - \bar{y})^2 &= \sum_h w_h^2 (\delta_h^{(\nu)})^2 + 2 \sum_{h \neq h'} w_h w_{h'} \delta_h^{(\nu)} \delta_{h'}^{(\nu)} \end{aligned}$$

WHERE $(-\nu)$ IS COMPLEMENT

$\sum_{\nu} \Rightarrow$ THE SECOND TERM GOES AWAY!

- TAKE A BALANCED SUBSET: $\sum_{\nu} \sum_{h \neq h'} \delta_h^{(\nu)} \delta_{h'}^{(\nu)} = 0$

$$\Rightarrow \sigma_{BRR}(\bar{y}) = \sigma(\bar{y}) !$$

JACKKNIFE:

G GROUPS

$$t^{(\nu)} = \frac{G}{G-1} \sum_{j \in \nu} \frac{y_j}{n_j} ; t(\bar{y}) = \frac{1}{G} \sum_{\nu=1}^G t^{(\nu)}$$

$$t^{\nu} = \sum_{j \in \text{GROUP } \nu} \frac{y_j}{n_j} \Rightarrow$$

$$G t - (G-1) t^{(\nu)} = t^{\nu}$$

$$-(G-1) (t^{(\nu)} - t) = t^{\nu} - t$$

$$\sigma(t) = \frac{1}{G(G-1)} \sum_{\nu} (t^{\nu} - t)^2 = \frac{G-1}{G} \sum_{\nu} (t^{(\nu)} - t)^2$$

SURV : DESIGN vs MODEL-BASED

Note Title

4/1/2007

MAIN REF: BINDER & ROBERTS (2003)

SEE ALSO: THOMPSON (1997, Ch. 5); HANSON, MADON & PEPPING (1983);
RUBIN-BLEUER & SCHIOPU-KRATINA (2005)

PRO-MODEL ARGUMENT:

IF THE DATA $\sim f(x, \theta)$, AND
THE MODEL IS TRUE, THEN THE
MODEL-BASED ESTIMATES
AND MODEL-BASED VARIANCES
ARE OPTIMAL!

STOCHASTIC COMPT:

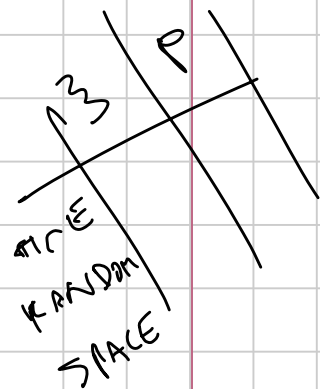
RESOLVED AT THE TIME OF
MEASUREMENT, SPACE = \mathbb{R}^n , RANDOM
VARS = (y_j, x_j)

PRO-DESIGN ARGUMENT:

WE HAVE TO HYPOTHESIZE ABOUT
A POTENTIAL ∞ POPULATION TO
JUSTIFY MODEL-BASED ARGUMENT,
WHILE IN THE DESIGN-BASED
APPROACH, THE FINITE POP~ PARAMETERS
ARE WELL-DEFINED, AND THERE IS
A FINITE POP~ SENSE OF CONSISTENCY.

STOCHASTIC COMPT:

RANDOM~ IN SAMPLING RESOLVED AT
TIME OF TAKING THE SAMPLE; THE VALUES
OF VARS OF INTEREST ARE CONSIDERED
TO BE FIXED; SPACE = $\{ \text{SAMPLES} \}$, $P(S)$
RANDOM VARS = I_{js}



SIMPLEST CASE:

ESTIMATING THE MEAN

MODEL:

Y_1, \dots, Y_n ARE SAID TO BE i.i.d.

$$E Y = \beta, \quad V Y = \sigma^2$$

\Rightarrow MODEL-BASED ESTIMATOR IS

$$E_{\beta} \bar{Y} = \beta; \quad V_{\beta} \bar{Y} = \sigma^2/n$$

DESIGN IS ASSUMED TO BE **IGNORABLE**:

$p(s)$ DOES NOT DEPEND ON Y_i 'S

DESIGN:

$Y_{j1}, \dots, Y_{jn} \leftarrow \text{POP} \sim \mathcal{U} = \{Y_1, \dots, Y_N\}$
POP ~ QUANTITIES; $b = \bar{y}_U = \frac{1}{N} \sum_{j=1}^N y_j$

SAMPLE ESTIMATES: SRS $\Rightarrow \hat{b} = \frac{1}{n} \sum_{j \in S} y_j$

$$E_p \hat{b} = b; \quad V_p \hat{b} = \frac{1}{n} (1-f) S^2$$

NON-SRS $\Rightarrow \hat{b} = \hat{b}_{HT} = \frac{1}{N} \sum_{j \in S} y_j / \pi_j$

REVIEW OF o , O , o_p , O_p

COMPARISON STRATEGY:

- TAKE $\{\text{DESIGN, MODEL}\}$ -BASED ESTIMATOR
- ESTABLISH $\{\text{MODEL, DESIGN}\}$ PROPERTIES
- GOOD ESTIMATORS: THOSE THAT SATISFY THE OTHER STRATEGY CONSISTENCY (OR EVEN EFFICIENCY)
- FRAMEWORK: FINITE POP ~ ASYMPTOTICS
 $N \rightarrow \infty, n \rightarrow \infty$: CLT IS APPLICABLE
- DISTINGUISH $o_p(n^q)$ UNDER DESIGN PROPERTIES;
 $O_p(n^q)$ UNDER MODEL PROPERTIES,
 $o(n^q)$ IN ASYMPTOTIC FRAMEWORK

PARAMETERS OF INTEREST:

$$\rightarrow \mu_j = E_3 Y_j \quad \text{@ INDIV LEVEL}$$

$$\rightarrow \bar{\mu} = \sum_{j=1}^N I_{js} \mu_j / n \quad \text{@ SAMPLE LEVEL}$$

(MOST OFTEN, $\mu_j = \mu \forall j$, ALTHOUGH NOT NECESSARILY)

$$\rightarrow \bar{\mu} = \sum_{j=1}^N \mu_j / N \quad \text{@ POP~ LEVEL}$$

→ RELATION BTW THE TWO: UNDER H_0 ,

$$\beta = \bar{\mu} + o_p(1) \quad (o_p(1)?)$$

$$\rightarrow b = \bar{y}_u = \frac{1}{N} \sum_{j=1}^N y_j \quad \text{@ POP~ LEVEL}$$

TARGETS OF INFERENCE:

β ← MODEL-BASED

b ← DESIGN-BASED

MODEL-BASED ESTIMATOR:

$$\hat{\beta} = \sum_{j \in S} I_{js} c_{js} Y_j$$

c_{js} ARE CHOSEN S.T.:

$$\beta = \beta + o_p(1) \Rightarrow \hat{\beta} = \bar{\mu} + o_p(1)$$

DESIGN-BASED ESTIMATOR:

$$\hat{b} = \sum_{j \in S} I_{js} d_{js} Y_j$$

d_{js} ARE CHOSEN S.T.:

$$\forall j, E_p(I_{js} d_{js}) = \frac{1}{N} + o(N^{-1}),$$

$$b = \hat{b} + o_p(1)$$

ESTIMATOR OF b

$$E_P(\hat{b}) = \sum_{j=1}^N y_j E_P(I_{js} d_{js}) = \frac{1}{N} \sum_{j=1}^N y_j + o_P(1) = b + o_P(1)$$

$$E_P(\hat{\beta}) = \sum_{j=1}^N y_j E_P(I_{js} c_{js}) \stackrel{!}{=} b + o(1)$$

THE MODEL IS TRUE

$$\Rightarrow E_3 b = \frac{1}{N} \sum E_3 y_j = \bar{\mu}; \quad b = \bar{\mu} + o_3(1)$$

$$E_3(b - \hat{\beta}) = E_3 \left[\bar{\mu} + o(1) - (\beta + o_3(1)) \right] =$$
$$= E_3 \left[(\bar{\mu} - \beta) + o(1) \right] = o(1)$$

BOTH $\hat{\beta}$ AND b DIFFER FROM THEIR E_3 BY $o(1) \Rightarrow$

$$E_P(\hat{\beta}) - b = E_3(E_P(\hat{\beta}) - b) + o(1)$$
$$= E_3(E_P(\hat{\beta} - b)) + o(1)$$
$$= E_P(E_3(\hat{\beta} - b)) + o(1)$$
$$= E_P(o(1)) + o(1) = o(1)$$

$$E_3 \hat{b} = \sum_{j \in S} I_{js} c_{js} \mu_j = \bar{\mu} + o_P(1)$$

$$E_3 b = \frac{1}{N} \sum_{j=1}^N E y_j = \bar{\mu}$$

$$E_P E_3 \hat{\beta} = E_P E_3 \hat{b} = \bar{\mu} + o(1)$$

DESIGN-BASED & TOTAL VARIANCES

$$\text{MODEL} \rightarrow \text{POP} \sim \rightarrow \text{SAMPLE}$$

$$O_3(N^{-1/2}) \quad O_p(n^{-1/2})$$

$$n = o(N)$$

$$V_{p3}(\hat{b}) = ?$$

$$V_p(\hat{b}) = V_p\left(\sum_{j \in S} I_{js} d_{js} y_j\right) = \sum_j \sum_{j'} \Delta_{jj'} y_j y_{j'}$$

$$\Delta_{jj'} = \text{Cov}_p(I_{js} d_{js}, I_{j's} d_{j's})$$

$$V_{3p}(\hat{b}) = V_3 E_p(\hat{b}) + E_3 V_p(\hat{b}) =$$

$$= V_3 [b + O_3(N^{-1/2})] + E_3 \sum_{jj'} \Delta_{jj'} y_j y_{j'} + o(n^{-1})$$

↑
DUE TO FPC

$$= V_3 \left[\frac{1}{N} \sum_j y_j \right] + \sum_{jj'} \Delta_{jj'} E_3 y_j y_{j'} + o(n^{-1}) + O_3(N^{-1/2})!$$

$$= \frac{1}{N^2} \sum_{jj'} \sigma_{jj'} + \sum_{jj'} \Delta_{jj'} (\sigma_{jj'} + \mu_j \mu_{j'}) + o(n^{-1})$$

$$\sigma_{jj'} = \text{Cov}_3(y_j, y_{j'}) ; \quad \frac{1}{N^2} \sum_{jj'} \sigma_{jj'} = \bar{\sigma}_{..}$$

i.i.d. (AND SOME OTHER WELL MIXING CASES)

$$\Rightarrow \bar{\sigma}_{..} = O(N^{-1}) = V_3 E_p \hat{b}$$

$$\Rightarrow V_{3p}(\hat{b}) = \sum_{jj'} \Delta_{jj'} (\sigma_{jj'} + \mu_j \mu_{j'}) + O(N^{-1}) + o(n^{-1})$$

$$\approx E_3 V_p[\hat{b}]$$

$$V(\hat{b}) / V_p(\hat{b}) = 1 + o(1) \Rightarrow \sigma_p(\hat{b}) / V_{3p}(\hat{b}) = 1 + o(1)$$

DESIGN-CONSISTENCY \Rightarrow TOTAL CONSISTENCY

$$\begin{aligned} \text{MSE}_{P_3}(\hat{\beta}) &= E_3 E_P (\hat{\beta} - b)^2 = \\ &= E_3 \left\{ \text{BIAS}_P^2(\hat{\beta}) + V_P(\hat{\beta}) \right\} \\ \text{UNDER MODEL} &\nearrow = O(N^{-1}) + O(n^{-1}) \\ &\quad \text{USUALLY} \quad \text{USUALLY} \end{aligned}$$

SPECIFIC EXAMPLE:

STRATIFIED SAMPLING : 2 STRATA

$$E_3 Y_j = \mu$$

$$\hat{\beta} = \bar{Y} = \frac{1}{n} \sum_{j \in S} Y_j ; \quad \hat{b} = W_1 \bar{Y}_{1n} + W_2 \bar{Y}_{2n}$$

$$E_P \hat{\beta} = (n_1 \bar{Y}_{1n} + n_2 \bar{Y}_{2n}) / n \neq E_P \hat{b} = \frac{N_1 \bar{Y}_{1n} + N_2 \bar{Y}_{2n}}{N} = b$$

DESIGN-INCONSISTENCY OF $\hat{\beta}$!

$$\text{BIAS}_P(\hat{\beta}) = E_P \hat{\beta} - b = \underbrace{\left(\frac{n_1}{n} - \frac{N_1}{N} \right)}_{\alpha} (\bar{Y}_{1n} - \bar{Y}_{2n}) = \alpha (\bar{Y}_{1n} - \bar{Y}_{2n})$$

MODEL IS TRUE \Rightarrow

$$E_3 \bar{Y}_{1n} = E_3 \bar{Y}_{2n} = \mu ; \quad E_3 E_P \hat{\beta} = \mu ; \quad E_P(\hat{\beta}) = \mu + o(1)$$

$$\text{ALSO, } E_3 b = \mu, \quad b = \mu + o(1)$$

$$\Rightarrow E_P \hat{\beta} = b + o(1) \quad - \text{ASYMPTOTIC}$$

DESIGN-UNBIASEDNESS

NOW: IF WE USE $\hat{\beta}$ TO ESTIMATE b ,

$$\begin{aligned} \text{MSE}_P(\hat{\beta}) &= E_P (\hat{\beta} - b)^2 = \\ &= \text{BIAS}_P^2(\hat{\beta}) + V_P(\hat{\beta}) \\ &= \alpha^2 (\bar{Y}_{1n} - \bar{Y}_{2n})^2 + (n_1 S_{1n}^2 + n_2 S_{2n}^2) / n^2 + o(n^{-1}) \end{aligned}$$

MSE_p($\hat{\beta}$) vs $V_p(\hat{b})$?

$$\alpha^2 (\bar{y}_{1n} - \bar{y}_{2n})^2 + \frac{n_1 S_{1n}^2 + n_2 S_{2n}^2}{n^2} \quad \text{vs} \quad \frac{N_1^2 S_{1n}^2}{N^2 n_1} + \frac{N_2^2 S_{2n}^2}{N^2 n_2}$$

TAKE $E_{\mathbb{Z}}$ OF BOTH SIDES:

$$E_{\mathbb{Z}} \text{MSE}_p(\hat{\beta}) = \alpha^2 \sigma^2 \left(\frac{1}{N_1} + \frac{1}{N_2} \right) + \frac{\sigma^2}{\bar{n}} = \frac{\sigma^2}{\bar{n}} + O(N^{-1}) = \frac{\sigma^2}{\bar{n}} + o(\bar{n}^{-1})$$

vs

$$E_{\mathbb{Z}} V_p(\hat{b}) = \sigma^2 \left(\frac{N_1^2}{N^2 n_1} + \frac{N_2^2}{N^2 n_2} \right) \geq \frac{\sigma^2}{\bar{n}}$$

WITH STRICT INEQUALITY WHEN
 $n_1/N_1 \neq n_2/N_2$

NON-LINEAR STATISTICS

MODEL-BASED ANALYSIS:

$$\Theta = g(\mu_1, \dots, \mu_m), \quad \hat{\Theta} = g(\bar{y}_{1n}, \dots, \bar{y}_{mn})$$

$$\begin{aligned} \Rightarrow \hat{\Theta} - \Theta &= \sum_{k=1}^m \frac{\partial g}{\partial \mu_k} (\bar{y}_{ks} - \mu_k) + o_{\mathbb{Z}}(n^{-1/2}) \quad \leftarrow \text{TAYLOR SERIES MODEL EXPN} \\ &= \sum_{j \in S} u_j / n + o_{\mathbb{Z}}(n^{-1/2}) = \sum_{j=1}^N \mathbb{I}_{j \in S} u_j / n + o_{\mathbb{Z}}(n^{-1/2}) \end{aligned}$$

WHERE $u_j = u_j(\mu) = \nabla g^T \cdot (y_{js} - \mu)$, $E_{\mathbb{Z}} u_j = 0$

$$\begin{aligned} E_p [g(\bar{y})] &= g(\mu) + \sum_{j=1}^N \pi_j u_j(\mu) / n + o_{\mathbb{Z}}(n^{-1/2}) \\ E_{\mathbb{Z}p} [g(\bar{y})] &= g(\mu) + o(n^{-1/2}) \end{aligned}$$

DESIGN-BASED ANALYSIS

POP ~ VALUE $\bar{Y}_{1n}, \dots, \bar{Y}_{mn}$

SAMPLE ESTIMATOR $\bar{y}_{1d}, \dots, \bar{y}_{md}$

$$g(\bar{y}_d) - g(\bar{Y}_n) = \sum_{j=1}^N I_{js} d_{js} u_j(\bar{Y}_n) + o_p(n^{-1/2})$$

$$\begin{aligned} E_p [g(\bar{y}_d)] &= g(\bar{Y}_n) + \sum \left[\frac{1}{N} + o(N^{-1}) \right] u_j(\bar{Y}_n) + o(n^{-1/2}) \\ &= \bar{g}(\bar{Y}_n) + o(n^{-1/2}) \\ E_{3p} [g(\bar{y}_d)] &= \bar{g}(\mu) + o(n^{-1/2}) \end{aligned}$$

$$\begin{aligned} V_p [g(\bar{y}_d)] &= \sum_{jj} \Delta_{jj} u_j(\bar{Y}_n) u_j(\bar{Y}_n) + o(n^{-1}) \\ &\approx V_{3p} [g(\bar{y}_d)] \end{aligned}$$

MODEL-ONLY PROPERTIES

USUALLY, z_3 -PROPERTIES OF $\hat{\beta}$ ARE WELL-KNOWN
E.G. MLE $\hat{\beta} \Rightarrow$ ASYMPT NORMALITY, ASYMPT EFFICIENCY

$$E_{z_3} \hat{b} = ?$$

$$E_{z_3} \hat{b} = E_{z_3} \left[\sum I_{js} d_{js} Y_j \right] = \sum I_{js} d_{js} \mu_j$$

$$= E_p \left[E_{z_3}(\hat{b}) \right] + o_p(1) \leftarrow LLN_p$$

$$= E_{z_3} \left[E_p(\hat{b}) \right] + o(1)$$

$$= E_{z_3} [\bar{Y}_n] + o(1) = \sum \mu_j / N + o(1)$$

$$E_{z_3} \hat{b} \text{ \& } (*) \Rightarrow E_{z_3} \hat{b} = \beta + o(1) -$$

ASYMPT MODEL-UNBIASED!

$$\begin{aligned}
V_3 \hat{b} &= ? \\
&= V_3 \left(\sum_j I_{js} d_{js} y_j \right) = \sum_{jj'} I_{js} I_{j's} d_{js} d_{j's} \sigma_{jj'} \\
&= E_p \left(\quad \right) + o(n^{-1}) \quad \Leftarrow LLN_p \\
&= \sum_{jj'} \delta_{jj'} \sigma_{jj'} + o(n^{-1}) \\
&= V_{3p} \hat{b} + o(n^{-1}) \quad \text{PROVIDED } \sum \mu_j \mu_{j'} \delta_{jj'} = o(n^{-1}) \\
&= E_3 \underbrace{V_p \hat{b}}_{\text{THE LEADING TERM OF } V_{3p} \hat{b}} + o(n^{-1})
\end{aligned}$$

$$\sum \mu_j \mu_{j'} \delta_{jj'} = o(n^{-1}) \Leftarrow$$

- $\mu_j = \mu \forall j, \quad P_p[n(s) = n] = 1$ (FIXED SIZE)
 $\Rightarrow \sum \mu_j \mu_{j'} \delta_{jj'} = 0$
- $\mu_j = 0 \forall j \Rightarrow \sum \mu_j \mu_{j'} \delta_{jj'} = 0$
 \uparrow ESTIMATING EQUATIONS!

MODEL VIOLATION

$$\text{IS } \hat{\beta} \xrightarrow{P} \beta?$$

$$\text{IF IT IS, IS } \sigma_3(\hat{\beta}) / V_3(\hat{\beta}) \rightarrow 1?$$

CRITICAL ASSUMPTIONS:

- SAMPLE DESIGN IS IGNORABLE
- STRUCTURE FOR THE MODEL MEANS / COVARS

MODEL VIOLATIONS:

- INFO SAMPLING: $E[y_j | I_{js} = 1] \neq \mu_j = E_3[y_j]$
- WRONG STRUCTURE OF μ_j OR $\sigma_{jj'}$

SAMPLING IS NOT IGNOREABLE

⇒ WHAT IS A REASONABLE TARGET FOR MODEL INFERENCE?

LET'S STILL CONSIDER / ASSUME

$$\beta = E_{\mathcal{Z}} \bar{y}_u + o(1) = \sum_{j=1}^J \mu_j / N + o(1) \quad (*)$$

SEEN EARLIER:

$$1) E_{\mathcal{Z}} [\hat{\beta}] = \beta + o(1)$$

$$2) \text{Var}_{\mathcal{Z}} [\hat{\beta}] = \text{Var}_{\mathcal{Z}} [\bar{y}_u] + o(1) \text{ UNDER CERTAIN ALTHOUGH REASONABLE CONDITIONS}$$

- ESTIMATING EQNS APPROACH IS ESPECIALLY BENEFICIAL!

THEN $\text{Var}_{\mathcal{Z}} [\hat{\beta}] \approx \text{Var}_{\mathcal{Z}} [\bar{y}_u]$ EVEN IF THE DESIGN IS NON-IGNOREABLE, BUT THE MEAN STRUCTURE IS CAPTURED

- IF THE MEAN STRUCTURE IS NOT MODELLED CORRECTLY, HOWEVER, THEN $E_{\mathcal{Z}} u_j(\bar{y}_u)$ MAY NOT BE ZERO, AND THEN $\text{Var}_{\mathcal{Z}} [\hat{\beta}] > \text{Var}_{\mathcal{Z}} [\bar{y}_u]$

PROPERTIES OF $\hat{\beta}$?

$$\hat{\beta} = \sum_j I_{js} c_{js} y_j \quad \leftarrow \text{BLUE OR SOMETHING}$$

• NON-IGNOREABLE DESIGN:

$$\mathcal{U} = \{0, \dots, 0, 1, \dots, 1\} \quad ; \quad P_{\mathcal{Z}} [u_j = 1] = \mu, \quad \beta = E_{\mathcal{Z}} [\bar{y}_u] = \mu$$

$$\pi_j = \begin{cases} p_0, & y_j = 0 \\ p_1, & y_j = 1 \end{cases} \quad \rightarrow \text{INFORMATIVE}$$

$$\mu_{j1} = E_{\mathcal{Z}} [y_j | I_{js} = 1] = P[y_j = 1 | I_{js} = 1] = \frac{\mu p_1}{(1-\mu)p_0 + \mu p_1}$$

$$\text{BIAS} [\bar{y}_s] = \mu_{j1} - \mu = \frac{\mu(1-\mu)(p_1 - p_0)}{(1-\mu)p_0 + \mu p_1}$$

• WRONG STRUCTURE

$$\mathcal{Z}: \begin{matrix} y_1 \\ y_2 \end{matrix} \sim N \left(\begin{pmatrix} \mu \\ \mu \end{pmatrix}, \frac{\sigma^2}{n} \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right)$$

$$\Rightarrow \text{BLUE } \hat{\mu} = 0.8y_1 + 0.2y_2$$

$$\mu_1 \neq \mu_2, \beta = E_{\mathcal{Z}} \hat{y}_1 = \frac{\mu_1 + \mu_2}{2}$$

$$\Rightarrow \text{BIAS } \hat{\mu} = E_{\mathcal{Z}} [\hat{\mu} - \beta] = 0.3(\mu_1 - \mu_2)$$

| | ρ | \mathcal{Z} |
|---------------------------|-----------|----------------------|
| ESTIMATOR | \hat{b} | $\hat{\beta}$ |
| CONSISTENCY FOR b ? | ✓ | WHEN THE MODEL HOLDS |
| CONSISTENCY FOR β ? | ✓ | WHEN THE MODEL HOLDS |

VARIANCE ESTIMATOR $V_{\rho}[\hat{b}] \approx V_{\rho\mathcal{Z}}[\hat{\beta}] !!$

V_{ρ} MODEL IS TRUE \Rightarrow OFTEN $V_{\rho}[\hat{b}] < V_{\rho}[\hat{\beta}]$

$V_{\mathcal{Z}}$ CUMBERSOME ... $(\text{INFO MX})^{-1}$

MODEL IS VIOLATED?

$$E_{\mathcal{Z}\rho} \hat{b} = \beta + o(1)$$

$$V_{\mathcal{Z}\rho}[\hat{b}] \approx V_{\rho}[\hat{b}]$$

UNDER MILD COND

$\hat{\beta}$ IS BIASED FOR β
MEANS ARE CORRECT
 $\Rightarrow V_{\rho}[\hat{\beta}] \approx V_{\mathcal{Z}}[\hat{\beta}]$

V_{ρ} 'S IN THOSE CASES ARE (ASYMPT) EQUIVALENT TO $V_{\mathcal{Z}}$ (IMPLIED BY LITERATURE ON M-ESTIMATES AND MISSPECIFIED MODELS (HUBER 67, 74; WHITE 82))
ARGUABLY, MSE IS THE BEST CRITERION?

MISSING DATA

- MISSING BY DESIGN

- TWO-PHASE SAMPLING:

X IS COLLECTED ON EVERYBODY
AT TIME t_1

Y IS COLLECTED ON A SUBSAMPLE
AT TIME t_2

- ROTATION SAMPLES:

DESCRIBE CPS

- SHORT AND LONG FORMS

- ATTRITION IN LONGITUDINAL SURVEYS

- ONLY A FRACTION OF THE UNITS
IN WAVE t WILL BE SURVEYED
IN WAVE $t+1$

- UNIT NON-RESPONSE:

(PRACTICALLY) NO INFORMATION
CAN BE OBTAINED FROM THE UNIT

- RESPONSE RATES: AAPOR DEFINITIONS
(INCLUDES CONTACTS, ELIGIBLE UNITS,
COMPLETES, INCOMPLETES, ...)

SHOW OF HANDS:

WHAT DO YOU THINK A
GOOD RESPONSE RATE IS?

- ITEM NON-RESPONSE:

FOR UNIT j , SOME DATA ARE COLLECTED, WHILE OTHERS ARE MISSING

FRAMEWORK:

UNIT j

VARIABLES

X_{ij} , $i=1, \dots, p$ - OF INTEREST

$f(x, \theta)$ - MODEL OF INTEREST

$z_{ij} = \mathbb{1}\{X_{ij} \text{ IS OBSERVED}\}$

MISSING DATA MECHANISM

$$p_{ij} = P[z_{ij}=1 | X_{ij}, X_{-ij}, X_{ij}', z_{ij}, z_{ij}', \psi]$$

MISSING COMPLETELY @ RANDOM:

$$P[z_{ij} | \text{SNPP}] = \psi$$

MISSING @ RANDOM

$$P[z_{ij} | \text{SNPP}] = P[z_{ij} | X_{-ij}, X_{ij}', z_{ij}, z_{ij}', \psi]$$

NOT MISSING @ RANDOM / INFORMATIVELY MISSING

$$P[z_{ij} | \text{SNPP}] = P[z_{ij} | X_{ij}, \dots]$$

THE MAIN THEOREM OF MISSING DATA:

THE MISSING DATA MECHANISM IS IGNORABLE

$$\in \text{MAR} \quad \& \quad \psi \cap \theta = \emptyset$$

- FACTORIZE LIKELIHOOD!

MISSING BY DESIGN:

$z_j \propto 1$ + RANDOMIZATION ✓

UNIT NON-RESPONSE:

ASSUME $P[z_j | \text{DESIGN VARS}] = \psi$

$\Rightarrow P[\text{IN SAMPLE}] = P[\text{SELECTION}] P[\text{RESPONSE}]$

$\Rightarrow \text{USE WEIGHT} = (P[\text{IN SAMPLE}])^{-1}$

ITEM NON-RESPONSE \Rightarrow ???

- LIKELIHOOD METHODS

- IMPUTATION METHODS

- HOT DECK

- MULTIPLE IMPUTATION

WEIGHTS

KORN & GRAUBARD (1999) - CH 4
PFEFFERMAN (1993)

PURPOSES OF WEIGHTS IN STAT PROGRAMS

= FREQUENCY WEIGHTS: HOW MANY IDENTICAL UNITS HAVE BEEN OBSERVED - ONLY MAKES SENSE FOR SIMPLE CATEGORICAL DATA / CONTINGENCY TABLES

USE: REPRODUCE CONTINGENCY TABLES, HISTOGRAMS, ... % PLOTS

= VARIANCE-ADJUSTING WEIGHTS:

IF $V y_j \propto \sigma^2 h_j$, THEN THE LKND IS
 $\ln L(y_j; \mu_j, \sigma_j^2) \sim \frac{1}{2} \sum_j \frac{(x_j - \mu_j)^2}{\sigma_j^2} \sim \sum_j \frac{(x_j - \mu_j)^2}{h_j}$

SO ONE CAN RUN THE ESTIMATION PROCEDURE WITH WEIGHT h_j^{-1} .

OFTEN, THESE WEIGHTS ARISE WHEN THE DATA ARE AGGREGATED BASED ON

GROUPS OF DIFFERENT SIZES:

$V[y_j] = \sigma^2 / n_j$, SO $h_j \propto n_j^{-1}$

WHICH MAKES IT SOMEWHAT LIKE FREQ WEIGHTS

- SAMPLING / SURVEY / PROBABILITY WEIGHTS TO CORRECT FOR DIFFERENTIAL PROBABILITIES OF SELECTION

FREQ INTERPRETATION:

EACH UNIT REPRESENTS w_j UNITS OF POPULATION

SAMPLING WEIGHTS ←

(1) PROBABILITY OF SELECTION, NT WAY:

$$w_j^{\text{BASE}} = 1/n_j$$

(2) NON-RESPONSE ADJUSTMENTS:

IF OUT OF n_j UNITS IN THE DESIGNED SAMPLE, THE DATA WERE COLLECTED ON \tilde{n}_j UNITS ONLY, THEN ASSUMING MAR | CLUSTER INDICATOR,

$$\hat{P}[\text{RESPONSE}] = \tilde{n}_j/n_j$$

$$w_j^{\text{NR}} = 1/\hat{P}[\text{RESPONSE}] = n_j/\tilde{n}_j$$

(3) FRAME UNDER COVERAGE/POSTSTRATIF.

IF CERTAIN DESIGN & OUTCOME VARS z_j ARE KNOWN FOR ALL POP~ UNITS. SEX-AGE-RACE: AVAILABLE FROM CENSUS, BUT ARE NOT KNOWN BEFORE SAMPLING, SO CANNOT BE USED TO STRATIFY BEFOREHAND

$$\sum_{j \in S} w_j^{\text{B}} w_j^{\text{NR}} z_{jk} = t_w[z_k]$$

USUALLY 0/1 CELL INDICATORS

$$\Rightarrow w_j^{\text{PS}} = \frac{t_w[z_k]}{t_w[z_k]} \quad \text{IF } z_{jk} = 1$$

THEN THE RESULTING FINAL WEIGHT WOULD BE

$$w_j = w_j^{\text{BASE}} w_j^{\text{NR}} w_j^{\text{PS}}$$

THE USE OF WEIGHTS:

$u_j(\theta)$ IS ESTIMATING EQ ~ FOR θ

$$0 = \sum_{j=1}^N u_j(\theta) \text{ FOR POP} \sim$$

$$0 = \sum_{j \in S} w_j u_j(\hat{\theta}) \text{ FOR SAMPLE}$$

- LINEAR STATISTICS: EXPLICIT SOLUTIONS

- NL (GLM): ITERATIVE ALGORITHMS

INEFFICIENCY OF WEIGHTED ESTIMATION

$$\bar{y}_w = \frac{t_w(y)}{t_w(1)} = \frac{\sum_j w_j y_j}{\sum_j w_j}$$

$$V_3[\bar{y}_w] = V_3\left[\frac{\sum_j w_j y_j}{\sum_k w_k}\right] = \frac{\sum_j w_j^2 \sigma^2}{\left(\sum_j w_j\right)^2}$$

$$DEFF_w = \frac{V_3[\bar{y}_w]}{V_{SRS}[\bar{y}]} = \frac{\sum_j w_j^2 \sigma^2}{\left(\sum_j w_j\right)^2} / \frac{\sigma^2}{n} =$$

$$= \left(\frac{1}{n} \sum_j w_j^2\right) / \left(\frac{1}{n} \sum_j w_j\right)^2 = 1 + \frac{V[w]}{E[w]^2} = 1 + CV_w^2$$

COEFFICIENT OF VARIATION OF WEIGHTS

SEE TABLE 4.4.1 OF K&G/4

$$INEFF_0 = 1 - \frac{V_{SRS}[\bar{y}]}{V_3[\bar{y}_w]} = 1 - \frac{1}{DEFF}; DEFF = (1 - INEFF_0)^{-1}$$

RATHER COMMON PRACTICE:
TRUNCATE THE WEIGHTS THAT ARE TOO HIGH

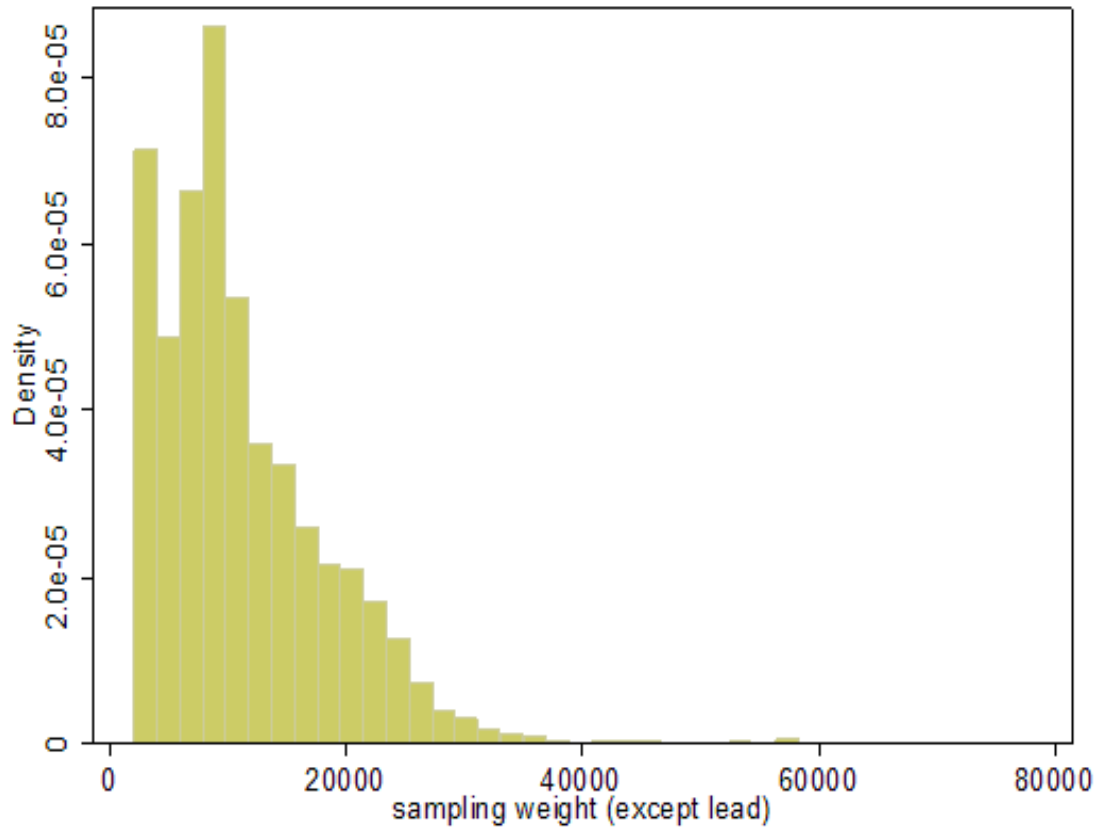
NHANES:

```
. use nhanes2
```

```
. sum finalw
```

| Variable | Obs | Mean | Std. Dev. | Min | Max |
|----------|-------|----------|-----------|------|-------|
| finalwgt | 10351 | 11318.47 | 7304.04 | 2000 | 79634 |

```
. di 1+r(Var)/( r(mean)*r(mean) )  
1.4164382
```



SERIOUS PFEFFERMANN (1993)

THE NEED FOR WEIGHTS ARISES WHEN THE DESIGNS ARE INFORMATIVE, AND MODEL IN THE POPULATION IS DIFFERENT FROM ONE IN THE DATA.

SIMPLE EXAMPLE:

DESIGN VARIABLE z_j : $P[j \in S] \propto z_j$

VARIABLE OF INTEREST y_j :

$\text{Cov}[y_j, z_j] > 0 \Rightarrow P[y_j > \mu_j | j \in S] > 1/2, E_S \bar{y}_S > \mu_y$

UNBIASED ESTIMATOR: REGRESSION

$$\bar{y}_R = \bar{y}_S + b(z_w - \bar{z}_S), \quad b = \sigma_{yz} / \sigma_z^2$$

REQUIRES THE KNOWLEDGE OF ALL VALUES z_j IN POP \sim TO MAKE DESIGN IGNORABLE!

NOTATION:

DESIGN VARS : $z_j \sim g(z, \varphi)$

SAMPLE DESIGN : $P(S) = P\{I_{jS}\} | Y, Z, \varphi$

MODEL OF INTEREST : $y_j \sim f(y | z, \theta)$

$Y = (Y_S, Y_{-S}) \leftarrow$ SAMPLED & NON-SAMPLED

C.F. MISSING DATA LITERATURE

JOINT DISTR \sim : $f(Y, z | \theta, \varphi) = f(Y | z, \theta) g(z, \varphi)$

SAMPLE DISTR \sim : OVER UNOBSERVED VARIABLES

$$f(Y_S, I_S, z | \theta, \varphi, Y) = \int f(Y_S, Y_{-S} | z, \theta) g(z, \varphi) P(I_S | Y_S, Y_{-S}, z, \varphi) dY_{-S}$$

NON-INFORMATIVE DESIGN:

$$f(Y_S, z | \theta, \varphi) = \int f(Y_S, Y_{-S} | z, \theta) g(z, \varphi) dY_{-S}$$

$$\text{E.G. } P(I_S | Y_S, Y_{-S}, z, \varphi) = P(I_S | z, \varphi)$$

WHEN z 'S ARE KNOWN FOR ALL POP \sim UNITS!

EXTENSIONS: 'IGNORABILITY FOR REGRESSION

$$f(Y_S | X_S, z, I) = f(Y_S | X_S, z)$$

$$\text{IF FURTHER } = f(Y_S | X_S) \rightarrow \text{OLS!}$$

IGNORING INFORMATIVE DESIGNS!

- BIAS OF POINT ESTIMATES
- POOR PERFORMANCE OF INFERENCE PROCEDURES
- NO INTERPRETATION OF $E_{p_3} \hat{\beta}_{OLS}$
- POPULATION-BASED CASE-CONTROL STUDIES:

$$Y = \begin{cases} 1, & \text{RARE} \\ 0, & \end{cases}$$

$$IP[I_{js}] = \begin{cases} 1, & Y_j = 1 \\ \pi_{c1}, & Y_j = 0 \end{cases}$$

MODEL: $f(Y|X, \theta)$ - LOGISTIC REGRESSION
(PSEUDO)MLE VS WEIGHTED?

IF THE MODEL HOLDS IN THE POPULATION,
 $\theta_{\text{INTERCEPT}}$ IS BIASED, BUT ALL θ_x 'S ARE IN
FACT MLE'S!

TESTING IGNORABILITY

REGRESSION CONTEXT:

vs. OPTIMAL / OLS / UNWEIGHTED / MLE $\hat{\beta}$
WEIGHTED / DESIGN-CONSISTENT $\hat{\beta}_w$

H_0 : DESIGN IS IGNORABLE $\Rightarrow \text{plim}_{n, N \rightarrow \infty} p_3 [\hat{\beta} - \hat{\beta}_w] = 0$

$$\lambda = \hat{D}^T [V(\hat{D})]^{-1} \hat{D}, \quad \hat{D} = \hat{\beta} - \hat{\beta}_w$$

IMPLEMENTATION:

$$Y_s = \underbrace{X_s}_{OLS} \hat{\beta} + \underbrace{W_s X_s}_{AUX \text{ REGRESSORS}} \hat{\gamma} + \text{ERROR} \quad \left. \vphantom{Y_s} \right\} \text{TEST}$$

$H_0: \gamma = 0$

GENERAL HAUSHMAN TEST:

$$\hat{\theta}_0, \hat{\theta}_1 \xrightarrow{p} \theta, \quad \hat{\theta}_0 \text{ IS ASYMPT EFFICIENT, } \sqrt{n}(\hat{\theta}_1 - \theta) \xrightarrow{d} N(0, V_1)$$

$$\Rightarrow \text{ASCOV}(\hat{\theta}_0, \hat{\theta}_1 - \hat{\theta}_0) = 0, \quad V(\hat{\theta}_1 - \hat{\theta}_0) = V(\hat{\theta}_1) - V(\hat{\theta}_0) \geq 0,$$

$$(\hat{\theta}_1 - \hat{\theta}_0)^T (V_1 - V_0)^{-1} (\hat{\theta}_1 - \hat{\theta}_0) \xrightarrow{d} \chi^2_{\dim \theta}$$

SIGNIFICANT COMPTS \Rightarrow IMPORTANT DESIGN VARS/INTERACTIONS!

$V[\theta_0], V[\theta_1] \leftarrow$ DESIGN-BASED ESTIMATORS

APPROACHES TO INCORPORATE WEIGHTS

- MODIFICATION OF MODEL-DEPENDENT

ESTIMATORS: REPLACE ALL SUMS

$\frac{1}{n} \sum_{i=1}^n \text{WHATEVER}$ BY $\frac{1}{N} \sum_{i=1}^n w_i \text{WHATEVER}$

$$b_{OLS} = \frac{\frac{1}{n} \sum x_i y_i - (\frac{1}{n} \sum x_i)(\frac{1}{n} \sum y_i)}{\frac{1}{n} \sum x_i^2 - (\frac{1}{n} \sum x_i)^2}$$

AMBIGUITY: MORE THAN A SINGLE
DESIGN-CONSISTENT ESTIMATOR MIGHT BE
AVAILABLE (E.G. b_{HT} , b_{RATIO} , b_{REG} , ...)

- MODELS THAT YIELD DESIGN-CONSISTENT ESTIMATORS

E.G. FIXED STRATUM EFFECT MODELS

$$Y_{hi} | \mu_h, \sigma_h^2 \sim N(\mu_h, \sigma_h^2), P(\mu_h, \ln \sigma_h) \propto 1$$

RANDOM EFFECT MODELS

$$\mu_h | \mu, \delta^2 \sim N(\mu, \delta^2), P(\mu, \ln \sigma_h, \ln \delta) \propto 1$$

- PSEUDO-LIKELIHOOD: MODIFY LIKELIHOOD SCORE EQNS

$$u(\theta) = \partial \ln f(y, \theta) \Rightarrow$$

$$\sum_{j \in S} w_j u_j(\theta) = 0$$

- ESTIMATING FUNCTIONS

$$\text{UNBIASED: } \mathbb{E}_Z g(y, \theta) = 0$$

$$\text{OPTIMAL: } g^* = \arg \min_g \mathbb{E}_Z (g^2) / \left[\mathbb{E}_Z \left. \frac{dg}{d\theta} \right|_{\text{TRUE } \theta} \right]^2$$

SAMPLE: DESIGN-UNBIASED

$$\mathbb{E}_P h(y_s, \theta) = g^*(y, \theta) \quad \forall \text{ POPN } y, \theta$$

$$\text{OPTIMAL: } h^* = \arg \min_h \mathbb{E}_P h^2(y_s, \theta) / \left[\mathbb{E}_P \frac{dh}{d\theta} \right]^2$$

h^* TURNS OUT TO BE H-T!

- WEIGHTS AS SURROGATE SUMMARIES OF DESIGN VARIABLES

$a = (a_1, \dots, a_N)$ IS AN ADEQUATE SUMMARY OF $z := P(I_j | z) = P(I_j | a)$
 $\Rightarrow (\pi_j)_{j=1}^N$ IS THE MINIMAL POSSIBLE SUMMARY

(TREATMENT EFFECT LITERATURE: PROPENSITY SCORE)

DON'T HAVE TO WORK OUT THAT WAY: (π_j) 'S MIGHT BE TOO COARSE!

- MLE FROM WEIGHTED DISTR ~:

MODEL $P(y_j, \alpha) = P(I_{js} | y_j)$ —
MODIFY DISTR ~ OF Y TO ACCOUNT FOR PROB ~ OF OBSERVING Y

$$f(y; \lambda) = \frac{P(y_j; \alpha) f(y_j; \theta)}{P[I_{js}=1]} =: f(y_j | I_{js}=1)$$

$$P[I_{js}=1] = \int P(y_j, \alpha) f(y_j; \theta) dy_j$$

$\Rightarrow \text{LKD}(\theta) \propto \prod f(y_j; \theta) / \left\{ \int P(y_j; \alpha) f(y_j; \theta) dy_j \right\}^n$
↑
PROPORTIONALITY INCLUDES α

EMPIRICAL WAY: ESTIMATE d , PLUG IN
JOINT WAY: MAX OVER d, θ

- EMPIRICAL LKLD:

$$\max \prod p_j \quad \left| \quad \sum p_j g(x_j) = 0 \right.$$

WHERE $\sum_{j=1}^N g(x_j) = 0$ IN POP

$$\theta = \int u(x) dF_N(x)$$

$\Rightarrow \hat{\theta} = \int u(x) dF_w(x), \quad F_w(x) = \sum p_j \delta(x_j)$

- CORRECTS FOR KNOWN POP ~ QUANTITIES

IMPUTATION:

A PROCEDURE OF COMING UP WITH A VALUE FOR MISSING DATA

- SINGLE IMPUTED VALUE:

$\{y_{ij}^*, \text{IMPUTATION INDICATOR } z_{ij}\}$
USUALLY, UNBIASED ESTIM~ IS FEASIBLE,
BUT VARIANCE ESTIM~ IS PROBLEMATIC

- MULTIPLE IMPUTED VALUE:

SEVERAL PLAUSIBLE VALUES y_{ij}^* -
EXTRA VARIABILITY PROVIDED AS NEEDED

IDEALLY y_{ij}^* COMES FROM PREDICTIVE DISTR~
 $f(y_{ij} | y_{(-i)j}, y_{i(-j)}, z_j; \theta, \alpha, Y)$
↑ ↑ ↑ ↑ ↑
OTHER VARS OTHER OBS DESIGN MODEL DESIGN
MISSING DATA ↓

- MEAN IMPUT~: $y_{ij}^* = \bar{y}_j$ (M.B. OVER IMPUT~ CELL)

- REG~ IMPUT~: $y_{ij}^* = X_{(-i)j}^T \delta$, $\delta \leftarrow$ (D.C.) REGRESSION
 $Y \rightarrow X_{(-i)}$, THE VARS
AVAILABLE FOR j TH CASE

- STOCH / REG~ IMPUT~:

$$y_{ij}^* = X_{(-i)j}^T \delta + \varepsilon_{ij}^*, \quad \varepsilon_{ij}^* \sim N(0, S^2) \text{ OR}$$

EMPIRICAL DISTR~ OF ε_{ij}

$$\text{OR } y_{ij}^* = \begin{cases} 0 \\ 1 \end{cases}, \quad P[1] = \Lambda(X_{(-i)j}^T \delta)$$

- HOT DECK IMPUT~:

$$y_{ij}^* = y_{i(-j)}, \text{ M.B. OVER IMPUT~ CELL}$$

(OFTEN STRATA OR PSU)

POPULAR WITH DATA PROVIDERS, BUT
PROPERTIES ARE NOT SO WELL KNOWN

- GOLD DECK INPUT~:

Y^*_{ij} = VALUE FROM EXTERNAL SOURCE
(REGISTER, CENSUS, LAST YEAR VALUE, ...)

RANDOM COMPONENTS:

} ← MODEL

P ← DESIGN

* ← INPUT~

} p* IF AVAILABLE?

VARIANCE ESTIM~ WITH IMPUTED DATA?

⚠ NON-RESPONSE BIAS ⇒ $MSE[\hat{\theta}] > V[\hat{\theta}]$?

- EXPLICIT FORMULAE?

$$V(t_{HT} | Y_{OBS}) = \sum_{i=1}^L \frac{(l t_i / \pi_i - t)^2}{l(l-1)}$$

← (i) t_i ARE UNBIASED FOR CLUSTER TOTALS

(ii) t_i 'S ARE \perp ← ADJUSTMENT/IMPUT~ CELLS

DO NOT CUT THROUGH CLUSTERS

- RESAMPLING:

BRR/BSMAR/JKNIFE → IMPUTE FOR THE ^(F) DATASET

→ ESTIMATE → RESULTING VARIABILITY

SUFFICES!

* IMPUT~ MUST PROVIDE CONSISTENT $\hat{\theta}^{*(F)}$!

* SOMETIMES MODIF~ MIGHT BE REQUIRED - E.G.

NO RESPONDENTS IN A PARTICULAR BSTRAP
SUBSAMPLE!

- MULTIPLE IMPUTATION

d -th IMPUTED DATA SET $\rightarrow \hat{\theta}_d, \psi_d; d=1, \dots, D$

\hookrightarrow COMBINE

- POINT ESTIMATE $\bar{\theta}_D = \frac{1}{D} \sum_{d=1}^D \hat{\theta}_d$

- WITHIN-IMPUTATION: $\bar{\psi}_D = \frac{1}{D} \sum_{d=1}^D \psi_d$

- BETWEEN-IMPUTATION: $B_D = \frac{1}{D-1} \sum_{d=1}^D (\hat{\theta}_d - \bar{\theta}_D)(\hat{\theta}_d - \bar{\theta}_D)^T$

- TOTAL VARIABILITY:

$$T_D = \bar{\psi}_D + \frac{D+1}{D} B_D$$

- NON-RESPONSE INFO:

$$\gamma_D = (1 + 1/D) B_D T_D^{-1}$$

- INFERENCE:

$$(\theta - \bar{\theta}_D) T_D^{-1/2} \sim t_D$$

$$V = (D-1) \left(1 + \frac{1}{D+1} \bar{W}_D B_D^{-1} \right)^2$$

\uparrow SATTERTHWAITTE APPROX

OR

$$V^* = (V^{-1} + V_{obs}^{-1})^{-1}$$

$$\hat{V}_{obs} = (1 - \hat{\gamma}_D) \left(\frac{V_{com} + 1}{V_{com} + 3} \right) V_{com}$$

$V_{com} = \text{DESIGN D.F.}$

IMPUTATIONS = ?

- TYPICAL ADVICE: 3-5

- EMPIRICAL EVIDENCE ON STABILITY?

- DESIGN DEGREES OF FREEDOM?

STAT 9100 SVY : DIFFICULT SITUATIONS

Note Title

4/25/2007

- ZERO SURVEY WEIGHTS? <http://www.stata.com/support/faqs/st>
- SUBDOMAIN ESTIMATION
 - ↳ SET WEIGHTS OF IRRELEVANT UNITS TO ZERO
 - ↳ D.F. = #PSUS WITH DOMAIN OBS - #STRATA WITH DOMAIN
- SINGLETON PSUS / STRATA (E.G. SELF-REPRESENTING, $\pi_j = 1$)
 - ↳ COLLAPSE STRATA: INCREASE D.F., OVERACCOUNT \downarrow BY STRAT
 - ↳ USE SELF-REPRESENTING PSUS AS PSEUDO-STRATA ← SSU
- LACK OF D.F.S (KORN & GRAUBARD 1995)
- PANEL DATA
 - OTHER EXAMPLES OF MISMATCH BETWEEN THE SURVEY POP & THE ANALYSIS POP:
BINDER & ROBERTS (2006)
- MATCHING & NESTED STRUCTURE - NEED π_{ij} ?
- SURVEYS OVER TIME
 - GOOD TIMES TO SURVEY? OPTIMIZE OVERLAP? SEASONALITY?
 - TIME IN SAMPLE EFFECT? ROTATING SURVEYS
- NON-SMOOTH ESTIMATORS (QUANTILES)
- VARIANCE ESTIM~ WITH IMPUTED DATA
- SURVEY BOOTSTRAP
 - SCALING; NON-RESPONSE, IMPUT~, POST-STRAT~ ADJ
- OUTLIERS & INFLUENTIAL OBSERVATIONS
(TRADITIONAL SENSE + BIG WEIGHTS)
- MULTIVARIATE METHODS
- MULTILEVEL MODELS
- NON-NORMAL ASYMPTOTICS?
 - χ^2 GOODNESS-OF-FIT → RAO-SCOTT MIXTURES
- NETWORK SAMPLING - EXTREMELY INFO!