

You should know...

- ▶ stratum, strata, stratified sampling
- ▶ why/when we prefer stratified sampling to SRS
- ▶ inference: use SRS theory in each stratum h
- ▶ $\hat{t}_{str} = \sum_{h=1}^H \hat{t}_h$, $\hat{V}(\hat{t}_{str}) = \sum_{h=1}^H (1 - n_h/N_h) N_h^2 s_h^2 / n_h$ (4.4)
- ▶ \bar{y}_{str} , $\hat{V}(\bar{y}_{str})$
- ▶ sampling weights:

$$\bar{y}_{str} = \sum_{h=1}^H \sum_{j \in S_h} w_{hj} y_{hj} / \sum_{h=1}^H \sum_{j \in S_h} w_{hj}$$
- ▶ proportional allocation: $w_{hj} = N/n$; self-weighting sample
- ▶ optimal allocation: $n_h \propto (N_h S_h) / \sqrt{c_h}$
- ▶ how to define strata; when is stratified sampling better (= smaller variance)
- ▶ **Final Exam: December 15, 9 - 11 a.m., EX 200 (255 McCaul Street)**
- ▶ **\$\$: Samuel Beatty Scholarship November 13**

... you should know...

- ▶ **HW:** Exercises 3.5, 3.13a, 3.15, 3.24; Examples 4.2, 4.3; Exercises 4.2, 4.12, 4.10; **Ex. 4.15 – new** ↗
- ▶ 3.5, 3.13a done on Friday; 3.15 posted online; 3.24 – you do it
- ▶ Example 4.2 done in class; Example 4.3 see text
- ▶ Exercise 4.2: posted online and see R code from Oct 23
- ▶ Exercise 4.12: optimal allocation for the agriculture data set

On p.97, we see we have estimates of S_h^2 for each of the 4 strata (for ACPRES92), and we know the population sizes N_h :

agrsr.dat

N_h	Stratum	Sample Size	s_h^2	optimal sample size
220	Northeast	21	7,647,472,708	69
1054	North Central	103	29,618,183,543	7
1382	South	135	53,587,487,856	122
422	West	41	396,185,950,266	101

agpop.dat 3087

300

77

p. 97

from agpop

- ▶ optimal: $n_h \propto N_h S_h^2 / \sqrt{c_h}$; if c_h 's all equal then $n_h \propto N_h S_h^2$; we use s_h^2 as estimates
- ▶ Exer. 4.10: It is **WRONG** on Rcode from October 23 – Friday

$$n_h \propto N_h S_h^2 / \sum (N_h S_h^2) \times 300$$

Exercise 4.10 (see R code from today)

	N_h	n_h	t_h
Bro	102	7	2+3+10+7
Phy	310	19	2+3+8+5+6+16
Soc	217	13	2+8+6
Hum	178	11	2+3
	<u>807</u>	<u>50</u>	

$$s_h^2 = \frac{2+9+50+49 - 22^2/7}{6}$$

$$\begin{aligned} & \sum_{i=1}^n (y_i - \bar{y})^2 \\ &= \sum y_i^2 - \frac{(\sum y_i)^2}{n} \end{aligned}$$

Cluster sampling Ch.5

stratified	cluster
variance within small strata $1, \dots, H$ population N_1, \dots, N_h observation y_{hj}	variance between small psu's $1, \dots, N$ – sampling unit ssu's M_1, \dots, M_n – observation unit observation y_{ij}

See Figure 5.1

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...) *500 such psu's*
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...)
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
- ▶ sample 5 suites at random
- ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
- ▶ 2 eggs in each nest chosen at random ...
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...) **500**
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...)
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
- ▶ sample 5 suites at random
- ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
- ▶ 2 eggs in each nest chosen at random ...
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...) *500 blk's*
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...) *400 total ss.*
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
- ▶ sample 5 suites at random
- ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
- ▶ 2 eggs in each nest chosen at random ...
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...)
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...)
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
- ▶ sample 5 suites at random
- ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
- ▶ 2 eggs in each nest chosen at random ...
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...)
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...)
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
- ▶ sample 5 suites at random
- ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
- ▶ 2 eggs in each nest chosen at random ...
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...)
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...)
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
- ▶ sample 5 suites at random
- ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
- ▶ 2 eggs in each nest chosen at random ...
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...)
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...)
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
 - ▶ sample 5 suites at random
 - ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
 - ▶ 2 eggs in each nest chosen at random ...
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...)
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...)
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
- ▶ sample 5 suites at random
- ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
- ▶ 2 eggs in each nest chosen at random ...
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...)
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...)
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
- ▶ sample 5 suites at random
- ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
- ▶ 2 eggs in each nest chosen at random ...
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...)
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...)
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
- ▶ sample 5 suites at random
- ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
- ▶ 2 eggs in each nest chosen at random ...
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...)
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...)
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
- ▶ sample 5 suites at random
- ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
- ▶ 2 eggs in each nest chosen at random ...
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling: Examples



- ▶ **Example, p.131:** 10,000 households; divide into blocks of 20 households (= ...)
- ▶ psu: sample 20 of the 500 blocks
- ▶ ssu: sample all 20 households on the block (total sample size = ...)
- ▶ cheaper, easier to implement
- ▶ values on a single block **more similar** than 20 values taken at random from all 10,000 households
- ▶ so less information than in an SRS of size 400
- ▶ **Example 5.2:** 400 students in a dorm, in suites of size 4
- ▶ sample 5 suites at random
- ▶ interview all 4 students
- ▶ **Example 5.6:** clutches (= nests) with ≥ 2 eggs each ...
- ▶ 2 eggs in each nest chosen at random ... *← n.b. 2 stage c.s.*
- ▶ **Example:** nearly all household surveys: <http://www.statcan.gc.ca/concepts/index-eng.htm>

... cluster sampling



- ▶ **one stage** cluster sampling: sample psu's by SRS, sample all ssu's
- ▶ **two stage** cluster sampling: sample psu's by SRS, sample ssu's using some probability method
- ▶ why?
- ▶ may not have a sampling frame of observation units (individuals in a city, customers of a store)
- ▶ population may be widely distributed
- ▶ population may occur in natural clusters
- ▶ example: nursing home residents
- ▶ usually, cost
- ▶ stratified sampling is **more efficient** (...)
- ▶ cluster sampling is **less efficient** (...)

... cluster sampling



- ▶ **one stage** cluster sampling: sample psu's by SRS, sample all ssu's
- ▶ **two stage** cluster sampling: sample psu's by SRS, sample ssu's using some probability method
- ▶ why?
- ▶ may not have a sampling frame of observation units (individuals in a city, customers of a store)
- ▶ population may be widely distributed
- ▶ population may occur in natural clusters
- ▶ example: nursing home residents
- ▶ usually, cost
- ▶ stratified sampling is **more efficient** (...)
- ▶ cluster sampling is **less efficient** (...)

... cluster sampling



- ▶ **one stage** cluster sampling: sample psu's by SRS, sample all ssu's
- ▶ **two stage** cluster sampling: sample psu's by SRS, sample ssu's using some probability method
- ▶ **why?**
 - ▶ may not have a sampling frame of observation units (individuals in a city, customers of a store)
 - ▶ population may be widely distributed
 - ▶ population may occur in natural clusters
 - ▶ example: nursing home residents
 - ▶ usually, cost
 - ▶ stratified sampling is **more efficient** (...)
 - ▶ cluster sampling is **less efficient** (...)

... cluster sampling



- ▶ **one stage** cluster sampling: sample psu's by SRS, sample all ssu's
- ▶ **two stage** cluster sampling: sample psu's by SRS, sample ssu's using some probability method
- ▶ why?
- ▶ may not have a sampling frame of observation units (individuals in a city, customers of a store)
- ▶ population may be widely distributed
- ▶ population may occur in natural clusters
- ▶ example: nursing home residents
- ▶ usually, cost
- ▶ stratified sampling is **more efficient** (...)
- ▶ cluster sampling is **less efficient** (...)

... cluster sampling



- ▶ **one stage** cluster sampling: sample psu's by SRS, sample all ssu's
- ▶ **two stage** cluster sampling: sample psu's by SRS, sample ssu's using some probability method
- ▶ why?
- ▶ may not have a sampling frame of observation units (individuals in a city, customers of a store)
- ▶ population may be widely distributed
- ▶ population may occur in natural clusters
- ▶ example: nursing home residents
- ▶ usually, cost
- ▶ stratified sampling is **more efficient** (...)
- ▶ cluster sampling is **less efficient** (...)

... cluster sampling



- ▶ **one stage** cluster sampling: sample psu's by SRS, sample all ssu's
- ▶ **two stage** cluster sampling: sample psu's by SRS, sample ssu's using some probability method
- ▶ why?
- ▶ may not have a sampling frame of observation units (individuals in a city, customers of a store)
- ▶ population may be widely distributed
- ▶ population may occur in natural clusters
- ▶ example: nursing home residents
- ▶ usually, cost
- ▶ stratified sampling is **more efficient** (...)
- ▶ cluster sampling is **less efficient** (...)

... cluster sampling



- ▶ **one stage** cluster sampling: sample psu's by SRS, sample all ssu's
- ▶ **two stage** cluster sampling: sample psu's by SRS, sample ssu's using some probability method
- ▶ why?
- ▶ may not have a sampling frame of observation units (individuals in a city, customers of a store)
- ▶ population may be widely distributed
- ▶ population may occur in natural clusters
- ▶ example: nursing home residents
- ▶ usually, cost
- ▶ stratified sampling is **more efficient** (...)
- ▶ cluster sampling is **less efficient** (...)

... cluster sampling



- ▶ **one stage** cluster sampling: sample psu's by SRS, sample all ssu's
- ▶ **two stage** cluster sampling: sample psu's by SRS, sample ssu's using some probability method
- ▶ why?
- ▶ may not have a sampling frame of observation units (individuals in a city, customers of a store)
- ▶ population may be widely distributed
- ▶ population may occur in natural clusters
- ▶ example: nursing home residents
- ▶ usually, cost
- ▶ stratified sampling is **more efficient** (...)
- ▶ cluster sampling is **less efficient** (...)

... cluster sampling



- ▶ **one stage** cluster sampling: sample psu's by SRS, sample all ssu's
- ▶ **two stage** cluster sampling: sample psu's by SRS, sample ssu's using some probability method
- ▶ why?
- ▶ may not have a sampling frame of observation units (individuals in a city, customers of a store)
- ▶ population may be widely distributed
- ▶ population may occur in natural clusters
- ▶ example: nursing home residents
- ▶ usually, cost
- ▶ stratified sampling is **more efficient** (...)
- ▶ cluster sampling is **less efficient** (...)

... cluster sampling



- ▶ **one stage** cluster sampling: sample psu's by SRS, sample all ssu's
- ▶ **two stage** cluster sampling: sample psu's by SRS, sample ssu's using some probability method
- ▶ why?
- ▶ may not have a sampling frame of observation units (individuals in a city, customers of a store)
- ▶ population may be widely distributed
- ▶ population may occur in natural clusters
- ▶ example: nursing home residents
- ▶ usually, cost
- ▶ stratified sampling is **more efficient** (...)
- ▶ cluster sampling is **less efficient** (...)

Mortality in Iraq

- ▶ “New study estimating number of dead in Iraq hotly contested” (Globe & Mail)
- ▶ “The human cost of the war in Iraq” (Economist)
- ▶ “A statistical study claims that many more Iraqis have died than was thought”
- ▶ “Mortality after the 2003 invasion of Iraq: a cross-sectional study” (The Lancet, 2006)
- ▶ “Iraqi death estimates called too high: methods faulted” (Science)

Mortality in Iraq

- ▶ “New study estimating number of dead in Iraq hotly contested” (Globe & Mail)
- ▶ “The human cost of the war in Iraq” (Economist)
- ▶ “A statistical study claims that many more Iraqis have died than was thought”
- ▶ “Mortality after the 2003 invasion of Iraq: a cross-sectional study” (The Lancet, 2006)
- ▶ “Iraqi death estimates called too high: methods faulted” (Science)

Mortality in Iraq

- ▶ “New study estimating number of dead in Iraq hotly contested” (Globe & Mail)
- ▶ “The human cost of the war in Iraq” (Economist)
- ▶ “A statistical study claims that many more Iraqis have died than was thought”
- ▶ “Mortality after the 2003 invasion of Iraq: a cross-sectional study” (The Lancet, 2006)
- ▶ “Iraqi death estimates called too high: methods faulted” (Science)

Mortality in Iraq

- ▶ “New study estimating number of dead in Iraq hotly contested” (Globe & Mail)
- ▶ “The human cost of the war in Iraq” (Economist)
- ▶ “A statistical study claims that many more Iraqis have died than was thought”
- ▶ “Mortality after the 2003 invasion of Iraq: a cross-sectional study” (The Lancet, 2006)
- ▶ “Iraqi death estimates called too high: methods faulted” (Science)

Mortality in Iraq

- ▶ “New study estimating number of dead in Iraq hotly contested” (Globe & Mail)
- ▶ “The human cost of the war in Iraq” (Economist)
- ▶ “A statistical study claims that many more Iraqis have died than was thought”
- ▶ “Mortality after the 2003 invasion of Iraq: a cross-sectional study” (The Lancet, 2006)
- ▶ “Iraqi death estimates called too high: methods faulted” (Science)

... mortality in Iraq

- ▶ Iraq body count: 48,693
- ▶ Burnham et al. : 601,027 (427,000 – 739,700)
- ▶ NEJM, Jan 2008 151,000 (104,000 – 223,000)
- ▶ based on IFHS study
http://www.emro.who.int/iraq/ifhs_faq.htm
- ▶ Journal of Peace Research, 2008: “Bias in epidemiological studies of conflict mortality”
- ▶ select a random main street
- ▶ choose a random cross street to the main street
- ▶ select a random household on the cross street to start the process
- ▶ interview that house and proceed to adjacent house until 40 houses have been surveyed

... mortality in Iraq

- ▶ Iraq body count: 48,693
- ▶ Burnham et al. : 601,027 (427,000 – 739,700)
- ▶ NEJM, Jan 2008 151,000 (104,000 – 223,000)
- ▶ based on IFHS study
http://www.emro.who.int/iraq/ifhs_faq.htm
- ▶ Journal of Peace Research, 2008: “Bias in epidemiological studies of conflict mortality”
- ▶ select a random main street
- ▶ choose a random cross street to the main street
- ▶ select a random household on the cross street to start the process
- ▶ interview that house and proceed to adjacent house until 40 houses have been surveyed

... mortality in Iraq

- ▶ Iraq body count: 48,693
- ▶ Burnham et al. : 601,027 (427,000 – 739,700)
- ▶ NEJM, Jan 2008 151,000 (104,000 – 223,000)
- ▶ based on IFHS study
http://www.emro.who.int/iraq/ifhs_faq.htm
- ▶ Journal of Peace Research, 2008: “Bias in epidemiological studies of conflict mortality”
- ▶ select a random main street
- ▶ choose a random cross street to the main street
- ▶ select a random household on the cross street to start the process
- ▶ interview that house and proceed to adjacent house until 40 houses have been surveyed

... mortality in Iraq

- ▶ Iraq body count: 48,693
- ▶ Burnham et al. : 601,027 (427,000 – 739,700)
- ▶ NEJM, Jan 2008 151,000 (104,000 – 223,000)
- ▶ based on IFHS study
http://www.emro.who.int/iraq/ifhs_faq.htm
- ▶ Journal of Peace Research, 2008: “Bias in epidemiological studies of conflict mortality”
- ▶ select a random main street
- ▶ choose a random cross street to the main street
- ▶ select a random household on the cross street to start the process
- ▶ interview that house and proceed to adjacent house until 40 houses have been surveyed

... mortality in Iraq

- ▶ Iraq body count: 48,693
- ▶ Burnham et al. : 601,027 (427,000 – 739,700)
- ▶ NEJM, Jan 2008 151,000 (104,000 – 223,000)
- ▶ based on IFHS study
http://www.emro.who.int/iraq/ifhs_faq.htm
- ▶ Journal of Peace Research, 2008: “Bias in epidemiological studies of conflict mortality”
- ▶ select a random main street
- ▶ choose a random cross street to the main street
- ▶ select a random household on the cross street to start the process
- ▶ interview that house and proceed to adjacent house until 40 houses have been surveyed

... mortality in Iraq

- ▶ Iraq body count: 48,693
- ▶ Burnham et al. : 601,027 (427,000 – 739,700)
- ▶ NEJM, Jan 2008 151,000 (104,000 – 223,000)
- ▶ based on IFHS study
http://www.emro.who.int/iraq/ifhs_faq.htm
- ▶ Journal of Peace Research, 2008: “Bias in epidemiological studies of conflict mortality”
- ▶ select a random main street
 - ▶ choose a random cross street to the main street
 - ▶ select a random household on the cross street to start the process
 - ▶ interview that house and proceed to adjacent house until 40 houses have been surveyed

... mortality in Iraq

- ▶ Iraq body count: 48,693
- ▶ Burnham et al. : 601,027 (427,000 – 739,700)
- ▶ NEJM, Jan 2008 151,000 (104,000 – 223,000)
- ▶ based on IFHS study
http://www.emro.who.int/iraq/ifhs_faq.htm
- ▶ Journal of Peace Research, 2008: “Bias in epidemiological studies of conflict mortality”
- ▶ select a random main street
- ▶ choose a random cross street to the main street
- ▶ select a random household on the cross street to start the process
- ▶ interview that house and proceed to adjacent house until 40 houses have been surveyed

... mortality in Iraq

- ▶ Iraq body count: 48,693
- ▶ Burnham et al. : 601,027 (427,000 – 739,700)
- ▶ NEJM, Jan 2008 151,000 (104,000 – 223,000)
- ▶ based on IFHS study
http://www.emro.who.int/iraq/ifhs_faq.htm
- ▶ Journal of Peace Research, 2008: “Bias in epidemiological studies of conflict mortality”
- ▶ select a random main street
- ▶ choose a random cross street to the main street
- ▶ select a random household on the cross street to start the process
- ▶ interview that house and proceed to adjacent house until 40 houses have been surveyed

... mortality in Iraq

- ▶ Iraq body count: 48,693
- ▶ Burnham et al. : 601,027 (427,000 – 739,700)
- ▶ NEJM, Jan 2008 151,000 (104,000 – 223,000)
- ▶ based on IFHS study
http://www.emro.who.int/iraq/ifhs_faq.htm
- ▶ Journal of Peace Research, 2008: “Bias in epidemiological studies of conflict mortality”
- ▶ select a random main street
- ▶ choose a random cross street to the main street
- ▶ select a random household on the cross street to start the process
- ▶ interview that house and proceed to adjacent house until 40 houses have been surveyed



Formulas

Notation

For cluster sampling

population quantities :

N # psu's \leftarrow sample units \checkmark

M_i # ssu \leftarrow observation units \checkmark

$K = \sum_{i=1}^N M_i$ total pop. size \checkmark

$t_i = \sum_{j=1}^{M_i} y_{ij}$ $t = \sum_{i=1}^N t_i$ \times

pop. quantities

$$\bar{y}_u = \sum_{i=1}^N \sum_{j=1}^{M_i} y_{ij} / K \quad S^2 = \sum_i \sum_j (y_{ij} - \bar{y}_u)^2 / (K-1)$$

$$\bar{y}_{iu} = \sum_{j=1}^{M_i} y_{ij} / M_i \quad S_i^2 = \sum_{j=1}^{M_i} (y_{ij} - \bar{y}_{iu})^2 / (M_i - 1)$$

unknown, we want to estimate it

sample n psu's m_i obsⁿ in each psu

\bar{y}_i sample mean in i th psu \leftarrow data

$$\hat{t}_i = M_i \bar{y}_i \quad \hat{t} = \sum_{i=1}^n \hat{t}_i \cdot \frac{N}{n} \leftarrow \text{est. of } t = \hat{t}_{\text{unb}}$$

$$s_t^2 = \sum_{i=1}^n (\hat{t}_i - \frac{\hat{t}}{n})^2 / (n-1)$$

\uparrow this is for Ch 5 + 6 ; but in Ch 5 $M_i \equiv M$

Two results: 1) $\hat{t} = \frac{N}{n} \sum \hat{t}_i$

$$SE(\hat{t}) = \sqrt{N^2 \left(1 - \frac{n}{N}\right) \frac{s_t^2}{n}} \quad (5.3)$$

2) $\hat{\bar{y}} = \frac{\hat{t}}{NM}$

$$SE(\hat{\bar{y}}) = \frac{1}{M} \sqrt{1 - \frac{n}{N} \frac{s_t^2}{n}} \quad (5.6)$$

$$\text{var}\left(\frac{\hat{t}}{NM}\right) = \frac{1}{(NM)^2} \cdot \text{var}(\hat{t})$$

Example 5.2

$N = 100$ dorm ~~rooms~~ suites per's
 $M = 4$ rooms in each suite
 $n = 5$ per's \leftarrow sample SRS
 $m = 4$

Cluster = suite

p 137

	1	2	3	4	5	
1	3.08	2.36				
2	2.60	.				
3	3.44	.				
4	3.04	.				
	<u>12.16</u>	<u>11.36</u>	<u>8.96</u>	<u>12.96</u>	<u>11.08</u>	$\leftarrow \hat{t}_i$'s

$$\hat{t} = \frac{N}{n} \sum \hat{t}_i = \frac{100}{5} (12.16 + \dots + 11.08) = 1130.4$$

$$\hat{\bar{y}} = \frac{\hat{t}}{NM} = \frac{1130.4}{400} = 2.826 \leftarrow \text{est. of } \bar{y}_u$$

2826
 11304
 4

2.826 OK??

$$s_t^2 = \frac{1}{4} \left\{ (12.16 - 1130.4/N)^2 + (11.36 - 1130.4/N)^2 + \dots + (11.08 - 1130.4/N)^2 \right\}$$

$\leftarrow 1130.4/N$ $= ? 2.256?$

$$SE(\hat{y}) = \sqrt{\left(1 - \frac{5}{100}\right) \frac{s_t^2}{n}} \cdot \frac{1}{M}$$

$$= 0.164 \quad n + bc.$$

$$2.826 \pm 1.96 \times 0.164 \quad (\text{Central limit theorem})$$

What's going on?

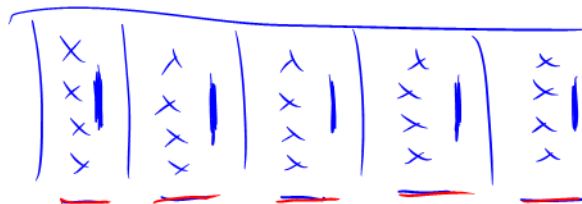
pop⁻ psu $i = 1, \dots, N$
 ssu $j = 1, \dots, M \leftarrow \text{ch. 5}$

pop⁻ values $y_{ij}; i = 1, \dots, N$
 $j = 1, \dots, M$

$$\sum_{i=1}^N \sum_{j=1}^M (y_{ij} - \bar{y}_u)^2 = \sum_i \sum_j (\underbrace{y_{ij} - \bar{y}_{iu}} + \underbrace{\bar{y}_{iu} - \bar{y}_u})^2$$

$$(k-1) S^2 = \sum_i \sum_j (y_{ij} - \bar{y}_{iu})^2 + \sum_i \sum_j (\bar{y}_{iu} - \bar{y}_u)^2 + 2 \times 0$$

$$\begin{array}{l} \text{total SS} \\ \text{in pop} \end{array} = \underbrace{\sum_i \sum_j (y_{ij} - \bar{y}_{iu})^2}_{\substack{\text{within SS} \\ \text{psu}}} + \underbrace{M \sum_i (\bar{y}_{iu} - \bar{y}_u)^2}_{\substack{\text{between SS} \\ \text{psu}}}$$



$$\begin{array}{l} \text{SSTO} \\ \uparrow \\ \text{SRS} \end{array} = \underbrace{\text{SSW}}_{\substack{\uparrow \\ \text{stratified}}} + \underbrace{\text{SSB}}_{\substack{\text{cluster sampling}}} - \left(\frac{\sum \hat{t}_h}{n} \right)^2$$

Def- ICC = "intraclass correlation coefficient"
or intraclass

$$= 1 - \frac{M}{M-1} \frac{\text{SS}_{\text{within}}}{\text{SS}_{\text{total}}}$$

It can be shown

$$\frac{V(\hat{t}_{clu})}{V(\hat{t}_{srs})} = \frac{MSB}{S^2} = \frac{SS_{betw}/(M-1)}{S^2} \approx 1 + (M-1)ICC$$