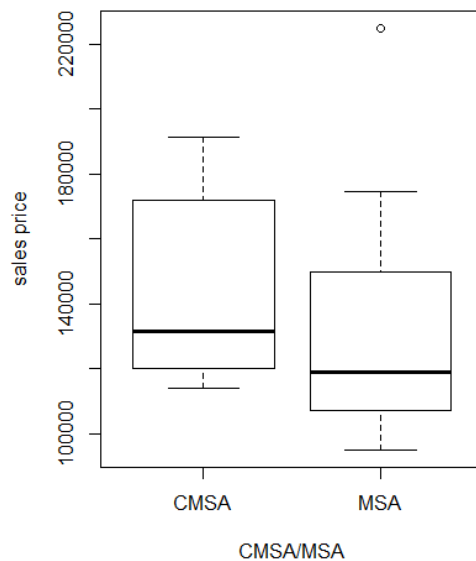**Kyle Beitz**

<u>**STAT 4155 Homework 2**</u>

**1. SMOG 5.4**

**a)**



One unusual feature i can see from this boxplot is the outlier in the MSA strata. This outlier will impact the accuracy of our estimates, so it should be accounted for in some way when making inferences. If the outlier seems to be an undeniable error in the data (say a sales price of -99999) it can be removed or corrected safely by speaking to the person who collected the data and fixing the error. If the outlier is not an error (in this case a sales price of 22500 seems possible), the outlier should not be removed as it is important to the model and should be accounted for. If still unsure about the outlier, using robust statistics may be the way to go. The mean and variance are very sensitive to outliers, whereas the median and interquartile range are less so. The latter should be considered over the former in situations where you are worried outliers may impact estimates greatly.

**b)** (on following page)

## Estimating means and totals from stratified samples

| | Strata Summary: | | | | | Estimate | SD(Est) | Confidence Intervals lower limit | upper limit | Relative margin of error |
|---|---|---|---|---|---|---|---|---|---|---|
| Stratum No. | Sample Size | Sample Size | Sample Mean | Sample SD | mean: | 131914.31 | 6323.74 | 119520.01 | 144308.61 | 9.40% |
| | | | | | total: | 35353034 | 1694762.26 | 32031361.01 | 38674706.99 | |
| 1 | 250 | 20 | 131045 | 31504.45 | | | | | | ? |
| 2 | 18 | 8 | 143988 | 28824.36 | | | | | | |

d.f. selector | infinite d.f. | confidence level 95% | z-multiplier 1.96

Ex. 5.2/5.3

Ex. 11.7a

2

Ex. 5.17

Ex. 11.7b

Number of decimals in answers

Case Study

Ex. 11.8

Clear out data    ?    Compare to SRS

| | Estimate | SD(Est) | Margin of error | Relative margin of error |
|---|---|---|---|---|
| mean: | 131914.31 | 19131.4900 | 37497.03 | 28.43% |
| | d.f. | 27 | t-multiplier | 1.96 |

totals: 268 28

**Comparing Two Stratum Means**

Strata to Compare | Estimate | SD(Est) | Confidence Intervals lower limit | upper limit

First
Second | | | d.f. | z-multiplier

Example 5.4

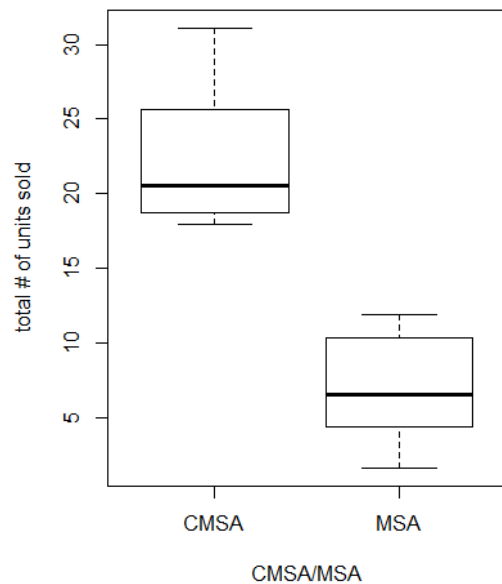Mean Typical Sales Price Per House For All Metropolitan Areas of the U.S.:

$$\underline{y}_{st} = 131914.31$$

Bound For Error of Estimation:

$$\underline{y}_{st} \pm 1.96(6323.74) \quad \text{<- Used 1.96 instead of 2 to match tool}$$

$$\underline{y}_{st} \pm 12394.5304$$

**c)** (on following page)

There are no outliers or unusual features of this boxplot to be concerned about.

**d)**



Total Number Of Houses Sold In All Metropolitan Areas of the U.S. In 1993:

$$N\underline{y}_{st} = 2151.4$$

Bound For Error of Estimation:

$N\underline{y}_{st} \pm 1.96(163.82)$     <- Used 1.96 instead of 2 to match tool

$$N\underline{y}_{st} \pm 321.0872$$

**e)**



Part b shows price. The gain in precision for stratifying is approximately 19% (RMoE of SRS - RMoE of Stratified Random Sampling). Part d shows total units sold. The gain in precision for stratifying is approximately 11%. The above screenshot shows square footage. The gain in precision for stratifying is approximately 26%. Thus, stratification produces the least gain in precision for total units sold. This is because total units sold does not differ greatly between the two strata (MSA and CMSA). The difference between the two strata's SDs is only about 1.5, while the difference between the SDs of the two strata for price is approximately 2700 and the difference between the SDs of the two strata for square footage is approximately 90. Stratification produces large gains in precision when the two strata greatly differ from each other, but within each strata the measurements are similar. This is why stratifying within total units sold does not create large gains in precision.

**f)**

**Estimating means and totals from stratified samples**

**Strata Summary:**

| No. | Stratum Sample Size | Sample Size | Sample Mean | Sample SD |
|-----|------|------|------|------|
| 1 | 250 | 20 | 131045 | 31504.45 |
| 2 | 18 | 8 | 143988 | 28824.36 |
| 3 | | | | |
| 4 | | | | |
| 5 | | | | |
| 6 | | | | |
| 7 | | | | |
| 8 | | | | |
| 9 | | | | |
| 10 | | | | |
| 11 | | | | |
| 12 | | | | |
| 13 | | | | |
| 14 | | | | |
| 15 | | | | |
| totals: | 268 | 28 | | |

| | Estimate | SD(Est) | Confidence Intervals lower limit | upper limit | Relative margin of error |
|---|---|---|---|---|---|
| mean: | 131914.31 | 6323.74 | 119520.01 | 144308.61 | 9.40% |
| total: | 35353034 | 1694762.26 | 32031361.01 | 38674706.99 | |

| d.f. selector | infinite d.f. | confidence level | z-multiplier |
|---|---|---|---|
| ◄ ► | | 95% | 1.96 |

◄ ►

Ex. 5.2/5.3

Ex. 5.17

Ex. 11.7a

Ex. 11.7b

Ex. 11.8

Case Study

Clear out data

?

**2** — Number of decimals in answers ◄ ►

?

| Compare to SRS | | | Relative |
|---|---|---|---|
| | Estimate | SD(Est) | Margin of error | Relative margin of error |
| mean: | 131914.31 | 19131.4900 | 37497.03 | 28.43% |
| | d.f. | 27 | t-multiplier | 1.96 |

**Comparing Two Stratum Means**

| Strata to Compare | Estimate | SD(Est) | Confidence Intervals lower limit | upper limit |
|---|---|---|---|---|
| First  2 | 12943.00 | 10166.3107 | -6982.6 | 32868.6 |
| Second  1 | | | d.f. | z-multiplier |
| | | | | 1.96 |

Example 5.4

Estimate The Difference In Average Typical Selling Price Between The Two Strata:

$$12943 \pm 1.96(10166.3107)$$
$$12943 \pm 19925.96897$$

CI: (-6982.6, 32868.6)

We cannot say that the houses in the CMSAs are, on average, higher priced than those in the MSAs because the above estimate and confidence interval includes 0. This means that there is a chance there is no real difference in average typical selling price between the two strata, even though the CMSAs sample mean is much larger.

## 2. SMOG 5.20

In this situation, we believe simple random sampling would work just as well because not much is changing between the months except for the month itself. There is no obvious reason why the batteries produced in one month would be much different than those produced in another month. The characteristics of the batteries produced month-to-month should be relatively homogenous (similar), so if simple random sampling is just as convenient, you might as well do that. If, however, there was reason to believe that batteries did vary greatly between months (say, for example, worker morale is higher in January than December so they therefore produce better working batteries in January), gains in precision would be made through stratifying by month.

## 3. SMOG 5.22

Stratification will produce large gains in precision over simple random sampling when measurements within each strata are homogenous. The more similar measurements are within each strata, the more payoff you get for using stratification (compared to simple random sampling). If each strata has a lot of variability within, stratification does not have as much of a positive impact on precision. To summarize, to gain the most precision with stratification, stratify the heterogeneous population into strata that have homogeneous measurements (measurements that are similar to each other). Putting all this into a more explicit framework, you gain the most precision using stratification when you have strata with wildly different means, small SD within each stratum, and large SD overall.

**4. SMOG 5.24**

If we used stratified random sampling we could make strata's through every 3 years of service and the firm and you could also make sub-strata within the above strata for male and female or job title. If we used simple random sampling we could do the sampling relatively easily and more cost-efficient. I would recommend stratified random sampling because it would give a more precise estimate because you are more likely to see a greater variance in simple random sampling between someone who is new and someone that has been with the firm for a long time. You could also see a greater variance between job titles as well.

**5. SMOG 5.32 (part a only)**

$$\underline{y}_{st} = \frac{1}{N} \sum_{i=1}^{L} N_i \underline{y}_i$$

= [(0.5*7.63)+(0.1*7.74)+(0.4*6.55)]

= 7.209

$$\hat{V}(\underline{y}_{st}) = \frac{1}{N^2} \sum_{i=i}^{L} N_i^2 (1 - \frac{n_i}{N_i})(\frac{S_i^2}{n_i})$$

$= \frac{1}{(1347+163+1095)^2} [(0.5)^2 * (\frac{.15^2}{1347}) + (0.1)^2 * (\frac{35^2}{163}) + (0.4)^2 * (\frac{0.11^2}{1095})]$ ←- Ignoring fpc because the sample sizes are large enough to justify not using them.

= 0.00000000517

$= \pm 2\sqrt{0.00000000517}$

= 0.00014

The mean time for the population of those giving anesthesia is 7.29 and the bound on error is 0.00014.

## 6. SMOG 5.35 (parts a, b, and e only)

|  | S:ACRES | W:ACRES | NC:ACRES | NE:ACRES |
|---|---|---|---|---|
|  | 251.71001 | 102.03 | 328.97 | 19.83 |
|  | 382.70999 | 473.92001 | 427.20999 | 65.99 |
|  | 45.610001 | 775.83002 | 366.92999 | 27.62 |
|  | 199.72 | 0.01 | 338.73001 | 46.61 |
|  | 138.42 | 1324.4 | 385.56 | 97.19 |
|  | 24.24 | 423.78 | 392.64001 | 11.64 |
|  | 83.07 | 177.33 | 162.24001 | 10.36 |
|  | 120.96 | 540.40997 | 479.89999 | 0.32 |
|  | 57.790001 | 751.52002 | 227.66 | 188.01 |
|  | 37.799999 | 1333.58 | 131.56 | 138.62 |
|  | 112.9 | 214.45 | 395.01999 | 35.34 |
|  | 19.709999 | 193.91 | 250.50999 | 32.53 |
|  | 130.88 | 208.16 | 300.97 | 76.99 |
|  | 5.9000001 | 598.69 | 134.2 | 234.39 |
|  | 558.31 | 1424.23 | 337.89001 | 85.11 |
|  | 236.77 | 631.38 | 799.60999 | 86.4 |
|  | 73.650002 | 669.52002 | 409.72 | 81.48 |
|  | 245.67999 | 1289.73 | 1481.5 | 89.04 |
|  | 37.549999 | 766.37 | 179.28 | 177.22 |
|  | 125.09 | 728.13 | 68.339996 | 121.91 |
|  | 32.970001 | 1465.79 | 1417.52 | 20.46 |
|  | 178.16 | 1868.33 | 8.7600002 | 17.71 |
|  |  |  |  |  |
|  |  |  |  |  |
| Standard Dev | 133.61 | 518.13 | 375.22 | 63.77 |
| Mean | 140.89 | 725.52 | 410.21 | 75.67 |
| N | 1376 | 418 | 1052 | 210 |
| n | 22 | 22 | 22 | 22 |

**Part A**

a) South Region:



**Estimation of means and totals for Simple Random Sampling, based on the sample mean**

| | |
|---|---|
| 22 | sample size |
| 140.891 | mean |
| 133.608 | SD |

251.71
382.71 Population
45.61 Size
199.72 484
138.42
24.24 fpc ?
83.07 0.955
120.96
57.79
37.8
112.9
19.71
130.88
5.9
558.31
236.77
73.65
245.68
37.55
125.09
32.97
178.16

Example 4.3

Example 8.13

Example 9.7

| | Estimate of Mean | SE | Margin of Error | Confidence Lower | Upper |
|---|---|---|---|---|---|
| | 140.891 | 27.83 | 57.876 | 83.015 | 198.767 |

| | Estimate of Total | SE | Margin of Error | Confidence Lower | Upper |
|---|---|---|---|---|---|
| | 68191.244 | 296338.086 | 616268.8 | 40179.3 | 96203.2 |

Relative Margin of Error  41.1%  ?

Confidence Level  95%

Number of digits in answers:  3

i)

ii) Using a 95% confidence interval we get a mean of 140.891 and a margin of error of 57.876.

b) West Region:



**Estimation of means and totals for Simple Random Sampling, based on the sample mean**

| | |
|---|---|
| 22 | sample size |
| 725.523 | mean |
| 518.128 | SD |

102.03
473.92 Population
775.83 Size
0.01 484
1324.4
423.78 fpc ?
177.33 0.955
540.41
751.52
1333.58
214.45
193.91
208.16
598.69
1424.23
631.38
669.52
1289.73
766.37
728.13
1465.79
1868.33

Example 4.3

Example 8.13

Example 9.7

| | Estimate of Mean | SE | Margin of Error | Confidence Lower | Upper |
|---|---|---|---|---|---|
| | 725.523 | 107.925 | 224.442 | 501.081 | 949.965 |

| | Estimate of Total | SE | Margin of Error | Confidence Lower | Upper |
|---|---|---|---|---|---|
| | 351153.132 | 1149190.617 | 2389872.7 | 242523.2 | 459783.1 |

Relative Margin of Error  30.9%  ?

Confidence Level  95%

Number of digits in answers:  3

i)

ii) Using a 95% confidence interval we get a mean of 725.523 and a margin of error at 224.442.

c) North Central Region:

i)

ii) Using a 95% confidence interval we get a mean of 410.215 and a margin of error at 162.536.

d) North East Region:



i)

ii) Using a 95% confidence interval we get a mean of 75.671 and a margin of error at 27.626.

Part B

To find the total region we take the sample size times the mean:
- South Region:
  - Total acreage: 193865.89
  - Bound on error: $2\sqrt{(1376)^2(1-\frac{22}{1376})\frac{(133.608)}{22}}=77762.4$
- West Region:
  - Total acreage: 303268.5
  - Bound on error: $2\sqrt{(418)^2(1-\frac{22}{418})\frac{(518.128)}{22}}=89885.9$
- North Central Region:
  - Total acreage: 431545.7
  - Bound on error: $2\sqrt{(1052)^2(1-\frac{22}{1052})\frac{(375.215)}{22}}=166542.6$
- North East Region:
  - Total acreage: 15890.99
  - Bound on error: $2\sqrt{(210)^2(1-\frac{22}{210})\frac{(63.774)^2}{22}}=5403.2$

Part E

| Estimating means and totals from Post-stratified samples | | | | | | | Relative margin of error |
|---|---|---|---|---|---|---|---|
| **Strata Summary** | | | | Estimate | SE | Confidence Intervals | |
| | | | | | | lower limit | upper limit | |
| Stratum No. | Stratum Size | Standard Deviation | Stratum Mean | mean | 309.08 | 33.26 | 243.89 | 374.27 | 21.09% |
| 1 | 210 | 63.77 | 75.67 | total | | | | | |
| 2 | 1052 | 375.22 | 410.21 | Sample Size | | confidence level | z-multiplier | ? |
| 3 | 418 | 518.13 | 725.52 | | | | | |
| 4 | 1376 | 133.61 | 140.89 | 88 | | 95% | 1.96 | 2 |

The mean acreage per county across the United States is 309.08.

| Stratum Summaries | | | | | | | | | off options |
|---|---|---|---|---|---|---|---|---|---|
| Stratum ID | Stratum Size | Stratum SD | Optimal Allocation | Equal Allocation | Mean Estimates | Stratum Costs | | Confidence level: | 95% |
| 1 | 1376 | 133.61 | 12.96 | 22 | | | | t: | 2.02 |
| 2 | 418 | 518.13 | 15.27 | 22 | | | | | |
| 3 | 1052 | 375.22 | 27.83 | 22 | | | estimated population mean: | | |
| 4 | 210 | 63.77 | 0.94 | 22 | | | | | |
| 5 | | | | | | | relative margin of error: | | |
| 6 | | | | | | For Mean | margin of error : | 67.5813 |
| 7 | | | | | | | standard error: | 33.48788779 |
| 8 | | | | | | For Total | margin of error : | 206528.4333 | Estimates are based on your allocation |
| 9 | | | | | | | standard error: | 102338.9851 | |

Comparing spreadsheet with hand calculations:

$$= \frac{1}{(3056)^2}\left[(1376)^2(1 - \frac{22}{1376})(\frac{(133.61)^2}{22}) + (418)^2(1 - \frac{22}{418})(\frac{(518.13)^2}{22}) + (1052)^2(1 - \frac{22}{1052})(\frac{(375.22)^2}{22}) + (210)^2(1 - \frac{22}{210})(\frac{(63.77)^2}{22})\right]$$

= 1121.4

$$= 2\sqrt{1121.4}$$

= 66.97

The mean acreage per county across the United States is 309.08 and the bound on error is 66.97 or 309.08 $\pm$ 66.97

## 7. SMOG 5.36



Using the allocation tool and plugging in all the known data, we can actually use this tool as a sample size calculator. Estimating mean acreage per county across the U.S., to get a margin of error of approximately 50,000 acres, we would use a sample size of 152. This gets us a margin of error of 49,924.5 acres (49.9245 in the tool because the data is in thousands).

The optimal allocation of these strata will be rounded to 44 for stratum 1, 2 for stratum 2, 52 for stratum 3, and 54 for stratum 4. The reason I rounded stratum 2 down to 2 instead of up is because I needed the four strata to add to 152 and the SD of stratum 2 is much lower than the other strata so it does not need as much attention as the other strata. To summarize, when estimating mean acreage, to get a margin of error of approximately 50,000 acres I would use a

sample size of 152, with stratum 1 (North Central) being of size 44, stratum 2 (North East) being of size 2, stratum 3 (South) being of size 52, and stratum 4 (West) being of size 54.