

## STAT 512 su 2021 hw 7

1. In a study of offspring sex-ratios in the Swedish-born population in [3], it is reported that among the 2,059,372 first-born children during the time period from 1932 until 2013, 1,058,701 were males, and among the 1,699,793 second-born children during the same time period, 874,369 were males.
  - (a) Build a 95% confidence interval for the difference  $p_1 - p_2$ , where  $p_1$  is the proportion of first-born children who are males and  $p_2$  is the proportion of second-born children who are males.

We will use the interval given by

$$\hat{p}_1 - \hat{p}_2 \pm z_{\alpha/2} \sqrt{\frac{\hat{p}_1(1 - \hat{p}_1)}{n_1} + \frac{\hat{p}_2(1 - \hat{p}_2)}{n_2}},$$

the upper and lower bounds of which we can compute in R with

```
alpha <- .05

n1 <- 2059371
n2 <- 1699793

p1.hat <- 1058701/n1
p2.hat <- 874369/n2

ME <- qnorm(1-alpha/2) * sqrt( p1.hat*(1-p1.hat)/n1 + p2.hat*(1-p2.hat)/n2)

L <- p1.hat - p2.hat - ME
U <- p1.hat - p2.hat + ME
```

The interval is  $(-0.001322975, 0.0007072851)$ .

- (b) Report the margin of error for the confidence interval.

The margin of error is 0.00101513.

- (c) Do you believe there is a difference between the proportion of males among first-born children and among second-born children? Base your answer on the data.

Since the confidence interval contains 0, we cannot conclude (at the 0.05 significance level) that there is any difference between the proportion of males among first-born and among second-born children.

2. In [2], continuing with the Swedish theme, a data set was analyzed in which the number of traffic accidents in a day on Swedish roads was measured when a speed limit was and was not enforced. The measurements were taken on several days during the years 1961 and 1962. Read in the data by installing the R package MASS and executing the following commands:

```
library(MASS)
X <- Traffic$y[Traffic$limit=="no"]
Y <- Traffic$y[Traffic$limit=="yes"]
```

- (a) Now construct a 95% confidence interval for the difference in the mean number of traffic accidents under enforcement and non-enforcement of a speed limit (non-enforcement minus enforcement). Assume non-Normality of the population distributions.

Since  $n_1 = 115$  and  $n_2 = 69$  we will use the interval

$$\bar{X} - \bar{Y} \pm z_{\alpha/2} \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}},$$

the upper and lower bounds of which we can compute in R with

```
alpha <- 0.05

n1 <- length(X)
n2 <- length(Y)

S1 <- sd(X)
S2 <- sd(Y)

X.bar <- mean(X)
Y.bar <- mean(Y)

L <- X.bar - Y.bar - qnorm(1-alpha/2) * sqrt( S1^2/n1 + S2^2/n2)
U <- X.bar - Y.bar + qnorm(1-alpha/2) * sqrt( S1^2/n1 + S2^2/n2)
```

We get the interval (1.785873, 6.64891).

- (b) Give an interpretation of your confidence interval. What is your recommendation about speed limits?

According to the confidence interval, there is a reduction in the average number of traffic accidents under enforcement of a speed limit by somewhere from approximately 1 to 7 accidents, with 95% confidence. There appears to have been some evidence in favor of enforcing speed limits (in Sweden in the 1960s).

3. The number of insects on the leaves of some plants were counted after the application of different pesticides. The data for pesticides “A” and “B”, taken from [1], can be read into R with:

```
data("InsectSprays")
X <- InsectSprays$count[InsectSprays$spray=="A"]
Y <- InsectSprays$count[InsectSprays$spray=="B"]
```

- (a) We know that these random samples were *not* drawn from Normal distributions. Explain why we know this without having to do any analysis.

Since the data are counts, they are discrete, but Normal random variables are continuous random variables.

- (b) Treating the random samples as though they *were* drawn from Normal populations, give a 95% confidence interval for the ratio  $\sigma_2^2/\sigma_1^2$ , where  $\sigma_2^2$  is the variance for pesticide “B” and  $\sigma_1^2$  is the variance for pesticide “A”.

We will use the interval

$$\left( \frac{S_2^2}{S_1^2} F_{n_1-1, n_2-1, 1-\alpha/2}, \frac{S_2^2}{S_1^2} F_{n_1-1, n_2-1, \alpha/2} \right),$$

the upper and lower bounds of which can be computed in R with

```
alpha <- 0.05

n1 <- length(X)
n2 <- length(Y)

S1 <- sd(X)
S2 <- sd(Y)

L <- S2^2/S1^2 * qf(alpha/2, n1-1, n2-1)
U <- S2^2/S1^2 * qf(1- alpha/2, n1-1, n2-1)
```

The interval is (0.2357854, 2.845125).

- (c) Comment on whether there is evidence to conclude that the two variances are unequal.

Since the 95% confidence interval for  $\sigma_2^2/\sigma_1^2$  contains the value 1, it is plausible, according to the interval, that the variances are equal. We therefore do not conclude (at the 0.05 significance level) that the variances are different.

- (d) Construct a 95% confidence interval for the difference  $\mu_1 - \mu_2$ , where  $\mu_1$  and  $\mu_2$  are the mean numbers of insects after the application of pesticides “A” and “B”, respectively.

Since we did not find strong evidence of a difference between the two variances, we will use the interval given by

$$\bar{X} - \bar{Y} \pm t_{n_1+n_2-2, \alpha/2} \sqrt{S_{\text{pooled}}^2 \left( \frac{1}{n_1} + \frac{1}{n_2} \right)},$$

the upper and lower bounds of which we can compute in R with

```
X.bar <- mean(X)
Y.bar <- mean(Y)
```

```
S.pooled <- sqrt( ( (n1 - 1) * S1^2 + (n2 - 1) * S2^2 ) / (n1 + n2 - 2) )
```

```
L <- X.bar - Y.bar - qt(1-alpha/2, n1 + n2 - 2) * S.pooled * sqrt( 1/n1 + 1/n2 )
U <- X.bar - Y.bar + qt(1-alpha/2, n1 + n2 - 2) * S.pooled * sqrt( 1/n1 + 1/n2 )
```

We get the interval  $(-4.643994, 2.977327)$ .

(e) Give an interpretation of the confidence interval.

Since the confidence interval for the difference in the two means contains zero, there is no evidence (at the 0.05 significance level) of a difference in the effectiveness of the two pesticides.

4. Let  $X_1, \dots, X_{n_1} \stackrel{\text{iid}}{\sim} \text{Normal}(\mu_1, \sigma^2)$  and  $Y_1, \dots, Y_{n_2} \stackrel{\text{iid}}{\sim} \text{Normal}(\mu_2, \sigma^2)$  be independent random samples (note that the population variances are both equal to  $\sigma^2$ ) with sample variances  $S_1^2$  and  $S_2^2$ , respectively.

(a) Show that

$$S_{\text{pooled}}^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n_1 + n_2 - 2}$$

is an unbiased estimator of  $\sigma^2$ .

We have

$$\begin{aligned} \mathbb{E}S_{\text{pooled}}^2 &= \frac{(n_1 - 1)\mathbb{E}S_1^2 + (n_2 - 1)\mathbb{E}S_2^2}{n_1 + n_2 - 2} \\ &= \frac{(n_1 - 1)\sigma^2 + (n_2 - 1)\sigma^2}{n_1 + n_2 - 2} \\ &= \sigma^2. \end{aligned}$$

(b) Use the fact that

$$\frac{(n_1 + n_2 - 2)S_{\text{pooled}}^2}{\sigma^2} \sim \chi_{n_1 + n_2 - 2}^2$$

to derive the upper and lower bounds of a  $(1 - \alpha)100\%$  confidence interval for  $\sigma^2$ .

In order to derive a  $(1 - \alpha)100\%$  confidence interval for  $\sigma^2$ , we put the pivot quantity into a probability statement as

$$P\left(\chi_{n_1 + n_2 - 2, 1 - \alpha/2}^2 < \frac{(n_1 + n_2 - 2)S_{\text{pooled}}^2}{\sigma^2} < \chi_{n_1 + n_2 - 2, \alpha/2}^2\right) = 1 - \alpha,$$

which we can rearrange to get

$$P\left(\frac{(n_1 + n_2 - 2)S_{\text{pooled}}^2}{\chi_{n_1+n_2-2, \alpha/2}^2} < \sigma^2 < \frac{(n_1 + n_2 - 2)S_{\text{pooled}}^2}{\chi_{n_1+n_2-2, 1-\alpha/2}^2}\right) = 1 - \alpha,$$

from which we can see that the interval

$$\left(\frac{(n_1 + n_2 - 2)S_{\text{pooled}}^2}{\chi_{n_1+n_2-2, \alpha/2}^2}, \frac{(n_1 + n_2 - 2)S_{\text{pooled}}^2}{\chi_{n_1+n_2-2, 1-\alpha/2}^2}\right)$$

is a  $(1 - \alpha)100\%$  confidence interval for  $\sigma^2$ .

5. Consider estimating a population proportion  $p$ .

- (a) What is the most conservative sample size (erring on the large size) required in order to build a 95% confidence interval for a population proportion  $p$  with a margin of error of at most 2%?

We find the smallest sample size  $n$  which satisfies

$$z_{0.05/2} \sqrt{1/2(1 - 1/2)/n} \leq 0.02,$$

which is  $n = \lceil (1.959964/0.02)^2 \cdot 1/2(1 - 1/2) \rceil = 2401$ .

- (b) What about with a margin of error of at most 1%?

We need  $n = \lceil (1.959964/0.01)^2 \cdot 1/2(1 - 1/2) \rceil = 9604$ .

- (c) What about with a margin of error of at most 0.5%?

We need  $n = \lceil (1.959964/0.005)^2 \cdot 1/2(1 - 1/2) \rceil = 38415$ .

- (d) If you quadruple the sample size, what happens to the width of the confidence interval?

Quadrupling the sample size halves the width of the confidence interval.

6. Researchers are interested in comparing the means  $\mu_1$  and  $\mu_2$  of two populations. A pilot study has suggested that the standard deviation  $\sigma_1$  of the first population is three times larger than the standard deviation  $\sigma_2$  of the second population.

- (a) The researchers have the resources to sample a total of 1,000 observations from the two populations. Find the number of observations  $n_1$  which should be drawn from population 1 and the number of observations  $n_2$  which should be drawn from population 2 such that the width of a confidence interval for  $\mu_1 - \mu_2$  is minimized.

According to our sample size formulas, we should take

$$n_1 = \left( \frac{\sigma_1}{\sigma_1 + \sigma_2} \right) n \quad \text{and} \quad n_2 = \left( \frac{\sigma_2}{\sigma_1 + \sigma_2} \right) n,$$

where  $n = 1000$ . If  $\sigma_1 = 3\sigma_2$ , then these formulas become

$$n_1 = \left( \frac{3}{4} \right) n \quad \text{and} \quad n_2 = \left( \frac{1}{4} \right) n,$$

so we should choose  $n_1 = 250$  and  $n_2 = 750$ . So we draw a larger sample from the population with the larger variance.

- (b) Suppose that the pilot study suggested  $\sigma_1 \approx 1$ . Recommend sample sizes  $n_1$  and  $n_2$  under which the margin of error of a 98% confidence interval for  $\mu_1 - \mu_2$  will be less than 0.10.

We should choose

$$n_1 = \left\lceil \left( \frac{3}{4} \right) n^* \right\rceil \quad \text{and} \quad n_2 = \left\lceil \left( \frac{1}{4} \right) n^* \right\rceil,$$

where  $n^* = (z_{.02/2}/0.10)^2(1 + 1/3)^2 = 962.1146$ . So we choose

$$n_1 = \lceil (3/4)962.1146 \rceil = 722 \quad \text{and} \quad n_2 = \lceil (1/4)962.1146 \rceil = 241.$$

Optional (do not turn in) problems for additional study from Wackerly, Mendenhall, Scheaffer, 7th Ed.:

- 8.61, 8.62, 8.64
- 8.70, 8.71, 8.74, 8.76
- 8.82, 8.83

## References

- [1] Geoffrey Beall. The transformation of data from entomological field experiments so that the analysis of variance becomes applicable. *Biometrika*, 32(3/4):243–262, 1942.
- [2] Åke Svensson. On a goodness-of-fit test for multiplicative poisson models. *The Annals of Statistics*, pages 697–704, 1981.
- [3] Brendan P Zietsch, Hasse Walum, Paul Lichtenstein, Karin JH Verweij, and Ralf Kuja-Halkola. No genetic contribution to variation in human offspring sex ratio: a total population study of 4.7 million births. *Proceedings of the Royal Society B*, 287(1921):20192849, 2020.