

## STAT 515 - Chapter 5 Supplement

Brian Habing - University of South Carolina

Last Updated: October 19, 2000

### S5 - Confidence Intervals for Variances

Just as we can use the t-distribution to form a confidence interval for the mean of a population, we can use the  $\chi^2$  and  $F$  distributions to make confidence intervals for the variance of a population, or the ratio of two variances. The logic is the same in both cases: use the sampling distribution to form a probability statement containing just the one unknown parameter, and then solve for that parameter.

#### S5.1 - The Confidence Interval for $\sigma$ or $\sigma^2$

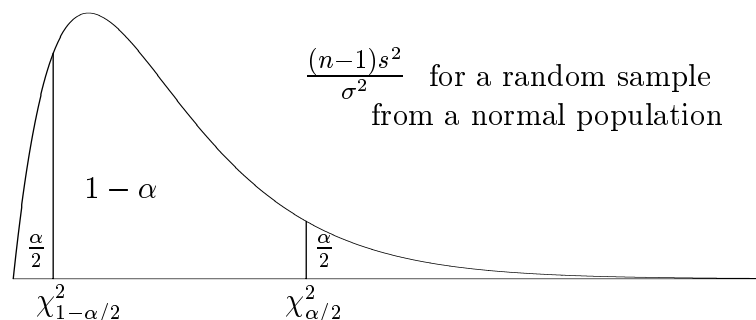
In section S4.2 we saw that if the random sample was drawn from a population that was normally distributed, then

$$\chi_{df=n-1}^2 = \frac{(n-1)s^2}{\sigma^2}$$

If we choose  $\chi_{\alpha/2}^2$  to be the value such that  $P(\chi_{df=n-1}^2 \geq \chi_{\alpha/2}^2) = \frac{\alpha}{2}$  and  $\chi_{1-\alpha/2}^2$  to be the value such that  $P(\chi_{df=n-1}^2 \leq \chi_{1-\alpha/2}^2) = \frac{\alpha}{2}$  we get the following:

$$P\left[ \chi_{1-\alpha/2}^2 \leq \frac{(n-1)s^2}{\sigma^2} \leq \chi_{\alpha/2}^2 \right] = 1 - \alpha \quad (S1)$$

This is illustrated in the figure below.



We can now solve the inequality in S1 for  $\sigma^2$  to get the confidence interval.

$$\begin{aligned}
& P \left[ \chi_{1-\alpha/2}^2 \leq \frac{(n-1)s^2}{\sigma^2} \leq \chi_{\alpha/2}^2 \right] \\
&= P \left[ \frac{1}{\chi_{1-\alpha/2}^2} \geq \frac{\sigma^2}{(n-1)s^2} \geq \frac{1}{\chi_{\alpha/2}^2} \right] \\
&= P \left[ \frac{1}{\chi_{\alpha/2}^2} \leq \frac{\sigma^2}{(n-1)s^2} \leq \frac{1}{\chi_{1-\alpha/2}^2} \right] \\
&= P \left[ \frac{(n-1)s^2}{\chi_{\alpha/2}^2} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi_{1-\alpha/2}^2} \right] = 1 - \alpha
\end{aligned}$$

The  $(1 - \alpha)100\%$  confidence interval for the population variance is thus:

$$\left( \frac{(n-1)s^2}{\chi_{\alpha/2}^2}, \frac{(n-1)s^2}{\chi_{1-\alpha/2}^2} \right) \tag{S2}$$

This can be changed to a confidence interval for the population standard deviation simply by taking the square root of both sides.

$$\left( \sqrt{\frac{(n-1)s^2}{\chi_{\alpha/2}^2}}, \sqrt{\frac{(n-1)s^2}{\chi_{1-\alpha/2}^2}} \right) \tag{S3}$$

Say we are interested in the standard deviation of a certain population. A sample of size  $n = 12$  ( $df = 11$ ) is gathered,  $s^2$  is found to be 20.2, the q-q plot for the data looks fairly normal, and it is desired to construct a 95% confidence interval. Using S3, we see that we need to determine  $\chi_{\alpha/2}^2$  and  $\chi_{1-\alpha/2}^2$ , where  $\alpha/2 = \frac{1-95\%}{2} = \frac{.05}{2} = 0.025$  and  $1 - \alpha/2 = 1 - 0.025 = 0.975$ . From Table XI we see that the values are  $\chi_{0.025}^2 = 21.9200$  and  $\chi_{0.975}^2 = 3.81575$ . Plugging these values into S3 gives:  $(\sqrt{10.137}, \sqrt{58.232}) = (3.18, 7.63)$ .

### S5.2 - The Confidence Interval for $\frac{\sigma_x^2}{\sigma_y^2}$

A confidence interval for the ratio of two variances could be constructed in the same way as one for a single variance. From S4.4 we see that if two samples  $x_1, x_2, \dots, x_{n_x}$  and  $y_1, y_2, \dots, y_{n_y}$  are drawn randomly from two populations that seem normal, then:

$$F_{df_x=n_x-1, df_y=n_y-1} = \frac{\frac{s_x^2}{s_y^2}}{\frac{\sigma_x^2}{\sigma_y^2}}$$

If we choose  $F_{\alpha/2}$  to be the value such that  $P(F_{df_x=n_x-1, df_y=n_y-1} \geq F_{\alpha/2}) = \frac{\alpha}{2}$  and  $F_{1-\alpha/2}$  to be the value such that  $P(F_{df_x=n_x-1, df_y=n_y-1} \leq F_{1-\alpha/2}) = \frac{\alpha}{2}$  we get the following:

$$P \left[ F_{1-\alpha/2} \leq \frac{\frac{s_x^2}{s_y^2}}{\frac{\sigma_x^2}{\sigma_y^2}} \leq F_{\alpha/2} \right] = 1 - \alpha \quad (S4)$$

Solving for  $\frac{\sigma_x^2}{\sigma_y^2}$  then gives us the  $(1 - \alpha)100\%$  confidence interval:

$$\left( \frac{\frac{s_x^2}{s_y^2}}{F_{\alpha/2}}, \frac{\frac{s_x^2}{s_y^2}}{F_{1-\alpha/2}} \right) \quad (S5)$$

It might be good practice to see if you can work out the steps between S4 and S5. Note that the text only gives the table for finding the  $F_{\alpha/2}$  values, and not the  $F_{1-\alpha/2}$  ones. However, SAS can be programmed to return those values and to form this confidence interval. Just remember that SAS gives the area in the lower end of the table, while the text (and the formulas above) use the upper end.

### S5.3 - Robustness... What if the Data Isn't Normal?

A statistical procedure is called robust if it performs well even when its assumptions aren't met. In the case of using the t-distribution to make inferences about the mean, the  $\chi^2$ -distribution for the variance, and the F-distribution for two variances,

we need to assume that the initial populations were normally distributed. The procedures that use the t-distribution are fairly robust however. That is, the procedures involving the use of the t-distribution to make inferences about  $\mu$  work fairly well even when the data isn't normally distributed. In general it will work well for small sample sizes ( $n \leq 30$ ) even if there are some doubts about the q-q plot. For large sample sizes it will work well for all but the worst q-q plots. The procedures discussed in S5.1 and S5.2 for making inferences about population variances are not robust at all. If there are any questions about the q-q plots, they should not be used. It is important to note that in Sections 8.3-8.4 and Chapter 9 we will see other uses of the  $\chi^2$  and  $F$  distributions that are robust. A particular distribution is never robust or non-robust, robustness is a property of the entire procedure that you are attempting.

In this course we are only covering some of the most common and basic methods for making inferences about a population. A variety of other methods are discussed in STAT 518 - Nonparametric Statistical Methods. A method is called nonparametric if it does not depend on the assumptions that the data follows a particular distribution (like the normal distribution). There are two reasons for not simply always using nonparametric methods. One is that they are somewhat more complicated to explain (see Section 7.5 for example). The second reason is that while the nonparametric tests are generally better when the assumptions are badly violated, the standard methods we are learning here are better when the assumptions are met.