

## 1. The F-test for Equality of Two Variances

Previously we've learned how to test whether two population means are equal, using data from two independent samples. We can also test whether two population variances are equal using sample data.

The  $F$  hypothesis test is defined as:

$$H_0 : \sigma_1^2 = \sigma_2^2$$

$$H_1 : \sigma_1^2 \neq \sigma_2^2$$

Test Statistic: 
$$F = \frac{S_1^2}{S_2^2} \quad \text{if } S_1^2 \geq S_2^2$$

$$F = \frac{S_2^2}{S_1^2} \quad \text{if } S_1^2 < S_2^2$$

where  $S_1^2$  and  $S_2^2$  are the sample variances. The more this ratio exceeds from 1, the stronger the evidence for unequal population variances.

The statistical significance of  $F$  is found by integrating an area of a cumulative  $F$  distribution.

As with the  $t$ - and  $\chi^2$  distributions, we have a different  $F$  distribution according to our degrees of freedom. However for the  $F$  statistic, we must consider the  $df$  associated with both variance estimates, i.e., ( $df_1, df_2$ ) degrees of freedom.

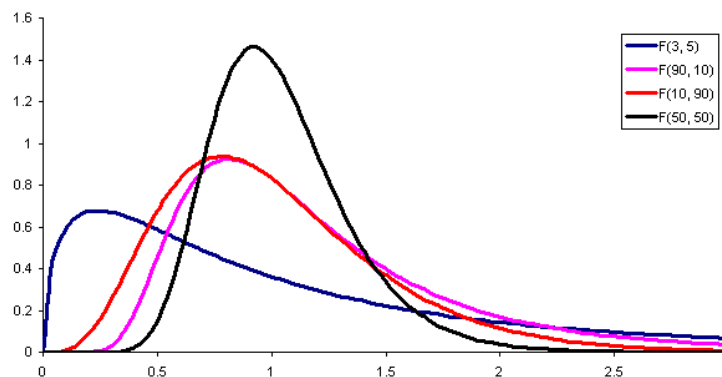
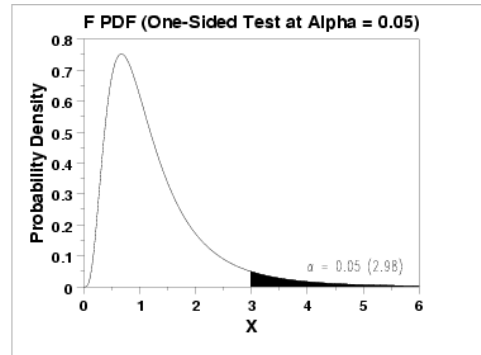
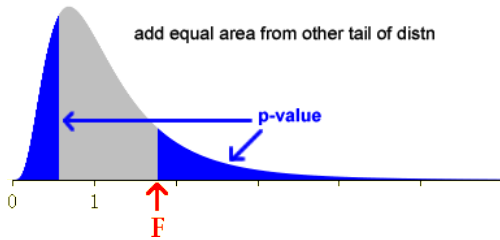


Figure 1. F Distributions for Various ( $df_1, df_2$ )



In theory, the F-distribution is two-tailed. That is, if  $F < 1$  we can integrate the distribution from 0 to  $F$ , or if  $F > 1$ , from  $F$  to  $+\infty$ . However (possibly due to considerations of numerical accuracy) the convention is to always integrate the positive tail.

For this reason we always place the larger of the two variance estimates in the numerator, and choose  $F$  distribution with:

$$(n_1 - 1, n_2 - 1) \text{ df when } s_1^2 \geq s_2^2, \text{ or}$$

$$(n_2 - 1, n_1 - 1) \text{ df when } s_1^2 < s_2^2.$$

### Excel

In Excel we can calculate  $F$  as the ratio sample variances, and then use the FDIST function to compute the  $p$ -value of  $F$ .

$$= \text{FDIST}(F, \text{df1}, \text{df})$$

Where:

- The larger variance should always be placed in the numerator.
- Multiply result  $\times 2$  for a two-tailed test.

To ensure the larger variance is always in the numerator, we use conditional logic in the cell formula for  $F$ :

Conditional logic in Excel: =IF(logical expression, formula if TRUE, formula if FALSE)

So, formula for  $F$ :

$$= \text{IF}(\text{var1} > \text{var2}, \text{var1}/\text{var2}, \text{var2}/\text{var1})$$

	A	B	C	D	E	F	G	H	I	J	K
1	F-test calculator										
2	Variable 1	Variable 2		Variance 1	2.9517						
3	1.48	7.55		Variance 2	3.5786						
4	1.75	3.75		n1	20						
5	0.78	0.1		n2	20						
6	2.85	1.1		df 1	19						
7	0.52	0.6		df 2	19						
8	1.6	0.52		F-test	1.2124	<-- =IF(var1>var2, var1/var2, var2/var1)					
9	4.15	3.3		p value (2-tail)	0.6789	<-- =2*IF(var1>var2,FDIST(F, df_1, df_2), FDIST(F, df_2, df_1))					
10	3.97	2.1									
11	1.48	0.58				where we have named the appropriate cells var1, var2,					
12	3.1	4.02				df_1, df_2, and F					
13	1.02	3.75									

**Assumption of F test of variances:** In the populations from which the samples were obtained the variable is normally distributed.

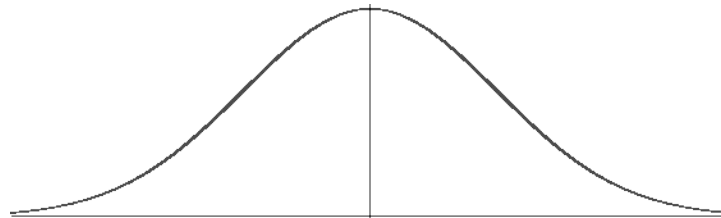
## 2. Credible and Confidence Intervals

You will recall that we've often made the distinction between credible intervals and confidence intervals. Credible intervals are derived using Bayesian statistics, and confidence intervals using classical statistical methods.

Also recall that the credible interval has a very clear and useful definition: it is the expected or distribution of some population parameter of, based on observed sample data. For example the 95% credible interval of a mean would tell us the estimated range for a population mean.

Again recall that the confidence interval, on the other hand, has a very convoluted and confusing definition, one not terribly helpful. In fact, the real reason people use confidence intervals is because they tend give them an interpretation that actually applies to the credible interval.

Fortunately, in several common cases, such as the z- and t-statistic, the confidence interval is, given mild assumptions, exactly equal to the credible interval. Here we will show why.



Suppose we have a random sample of  $n$  cases. We can compute the sample mean and sample standard deviation, and from these construct an expected sampling distribution. The mean of our sampling distribution would be our sample mean, and the standard deviation would be  $s/\sqrt{n}$ .

Suppose, for convenience, that our sample mean is exactly 0, and our standard error is exactly 1.

Now suppose we wish to estimate, given this sampling distribution, the probability (actually a probability density) of drawing a new sample of exactly  $n$  cases that has a sample mean of  $-1$ .

The usual formula for the probability density of a normal distribution is:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$

But because our mean is 0 and standard deviation is 1 this simplifies to:

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-.5z^2}$$

If  $z = 1$ , the above formula gives a value of 0.242.

Here's where the idea of "if I'm  $x$  units away from you, you're also  $x$  units away from me" applies. Suppose now that our population mean actually were  $-1$ , then our actual mean, 0, would correspond to a  $z$  of 1. And using the same formula above, the density of this value would also be 0.242.

Further, we could follow the same arrangement for every possible value for the population mean of our sampling distribution. In other words, in every case:

$$P(\bar{X} = c_1 | \mu_{\bar{X}} = c_2) = P(\mu_{\bar{X}} = c_1 | \bar{X} = c_2) = P(\mu_{\bar{X}} = c_2 | \bar{X} = c_1)$$

Where  $c_1$  and  $c_2$  are any two values.

The only requirement for this to work is the assumption that the standard error,  $s_{\bar{X}}$ , is the same for every possible value of  $\mu_{\bar{X}}$ . This is true for both the  $z$  and the  $t$  distributions.

The assumption of a constant standard error of the sampling distribution, however, is not true in the cases of a sample proportion. Recall that there the standard error is estimated as:

$$\sqrt{\frac{(p)(1-p)}{N}}$$

So as  $p$  changes, so does the standard error of  $p$ . Therefore we need another approach to construct the credible interval of a proportion.

### 3. Credible Interval for a Proportion

Class demonstration

