

Learning the F-Test: Comparing Variances in Python

Authored by
Mohammed loot

November 8, 2025

RECOMMENDED CITATION

Mohammed loot (2025). *Learning the F-Test: Comparing Variances in Python*. PSYCHOLOGICAL STATISTICS. Retrieved from <https://statistics.arabpsychology.com/?p=12683>

The Foundation: Understanding the F-Test for Variance Comparison

The **F-test**, named in tribute to the pioneering statistician Sir Ronald Fisher, is a cornerstone of classical statistics. Its fundamental purpose is to rigorously determine whether the underlying **population variances** of two independent data samples are statistically equivalent. This comparison is not merely academic; it is a critical prerequisite step in advanced data science and research pipelines, especially before conducting parametric tests like the two-sample t-test, which often relies on the crucial assumption of equal variances, also known as **homoscedasticity**.

In contrast to statistical tests that analyze differences in means (averages), the F-test zeroes in exclusively on the spread, dispersion, or variability within two distinct populations. The resulting F-statistic is derived simply as the ratio of the two observed sample variances. Logically, if the true population variances are identical, this ratio should approximate 1. A calculated F-statistic that deviates substantially from 1 suggests that the initial assumption of variance equality--the null hypothesis--must be challenged and potentially rejected.

Managing and verifying variability is paramount in building reliable statistical models. Datasets exhibiting a high degree of uncontrolled variability can undermine the validity of inferential conclusions. Consequently, confirming the equality of variances is a necessary procedural step in numerous inferential statistical methods. This comprehensive tutorial will guide you through a reliable and transparent methodology for performing this critical test directly within the powerful Python programming environment, leveraging standard scientific libraries.

Formulating the Hypotheses: Defining the F-Test Framework

Every procedure in **hypothesis testing** demands a precise definition of the competing statements: the null hypothesis (H_0) and the alternative hypothesis (H_1). These statements serve as the guiding principles for interpreting the calculated F-statistic and the subsequent **p-value**, which together determine our statistical conclusion.

The **null hypothesis** (H_0) is always the conservative default assumption, proposing that there is absolutely no statistically significant difference between the two population variances (σ_1^2 and σ_2^2). Conversely, the alternative hypothesis (H_1) asserts that a meaningful difference in variance exists. For the standard two-tailed F-test concerning the equality of variances, the formal mathematical statements are structured as follows:

H0: $\sigma_1^2 = \sigma_2^2$ (The population variances are equal, confirming homoscedasticity.)

H1: $\sigma_1^2 \neq \sigma_2^2$ (The population variances are *not* equal, indicating heteroscedasticity.)

Our objective throughout this process is to assess whether the empirical evidence provided by our sample data is strong enough to justify rejecting the conservative null hypothesis (H_0). If the calculated p-value falls below our predetermined threshold of significance (alpha, which is typically set at 0.05), we conclude that we have sufficient statistical evidence to reject H_0 and affirm that the spread or variability of the two populations is indeed fundamentally different.

Prerequisites and Data Preparation in the Python Environment

Effective execution of the F-test in Python necessitates leveraging the extensive capabilities of its scientific computing ecosystem. We rely heavily on the powerful numerical operations provided by the [NumPy](#) library for efficient array handling and precise variance calculation, alongside the statistical functions housed within [SciPy](#), particularly for accessing the F-distribution's cumulative density function (CDF). It is imperative that both of these libraries are correctly installed and imported into your working environment before any analysis can commence.

For the purpose of this practical demonstration, we will work with two hypothetical datasets. These samples represent measurements--perhaps performance scores--collected from two distinct treatment groups (Group X and Group Y). Our investigative question is whether the inherent variability in scores differs significantly between these two groups, regardless of their average performance.

We define our two data arrays, x and y , containing ten observations each ($n_1 = 10, n_2 = 10$):

$x =$

$y =$

A crucial underlying assumption of the F-test is that the data within each population is approximately [normally distributed](#), and that the samples were collected independently of one another. While the F-test is known to be quite sensitive to departures from normality, its core mathematical focus remains the accurate calculation of the ratio between the observed sample variances.

Developing the Custom F-Test Implementation in Python

A dedicated, centralized function for the two-sample F-test for variance equality is currently not standardly included within the primary SciPy statistics module, unlike functions for the t-test or ANOVA. To maintain control, transparency, and statistical rigor, we must construct a custom function. This function will seamlessly integrate [NumPy](#) for calculating the required variances and SciPy's F-distribution utilities for determining the corresponding p-value.

The F-statistic is mathematically defined as the ratio of the sample variances: $F = s_{\{1\}}^2 / s_{\{2\}}^2$. A critical statistical detail in computing this ratio is the proper calculation of the [sample variance](#). This requires applying [Bessel's correction](#), which involves dividing the sum of squared deviations by $n-1$ (the degrees of freedom) rather than n (the sample size). In NumPy's variance calculation, this correction is activated by setting the `ddof` (Delta Degrees of Freedom) parameter to 1.

We define the function `f_test(x, y)`, which takes our two data arrays as input and reliably returns both the calculated F-statistic and the associated p-value. The code below illustrates the calculation, emphasizing the use of `ddof=1` for accurate sample variance and the definition of the numerator and denominator [degrees of freedom](#) (dfn and dfd, respectively), which are essential inputs for the F-distribution functions.

```
import numpy as np
import scipy.stats

#define F-test function
def f_test(x, y):
    x = np.array(x)
    y = np.array(y)
    f = np.var(x, ddof=1)/np.var(y, ddof=1) #calculate F test statistic
    dfn = x.size-1 #define degrees of freedom numerator
    dfd = y.size-1 #define degrees of freedom denominator
    p = 1-scipy.stats.f.cdf(f, dfn, dfd) #find p-value of F test statistic
    return f, p

#perform F-test
f_test(x, y)

(4.38712, 0.019127)
```

The calculation of the p-value using `1 - scipy.stats.f.cdf(f, dfn, dfd)` is specific to the F-test structure where we conventionally place the larger [sample variance](#) in the numerator, making the test inherently one-tailed in its calculation. The cumulative distribution function (CDF) returns the probability of observing an F-statistic *less than or equal to* our calculated value. By subtracting this result from 1, we isolate the probability residing in the upper tail of the [F-distribution](#)--which represents the probability of observing an F-statistic as extreme, or more extreme, than the one calculated, thus yielding the precise p-value for our test.

Interpreting the Output and Drawing Statistical Conclusions

Upon executing the custom function with our sample data sets (x and y), the output provides the two essential metrics required for statistical inference: the calculated F test statistic and its corresponding p-value. These empirical findings form the basis for making a decisive judgment about the equality of the population variances.

From the output, we observe that the F test statistic is calculated as **4.38712**, and the corresponding p-value is **0.019127**. To properly interpret these figures, we must compare the **p-value** against our predetermined level of significance (α). Standard practice dictates setting α at 0.05, signifying our maximum acceptable risk--a 5% chance--of committing a **Type I error** (falsely rejecting a true null hypothesis).

Given that the calculated p-value of 0.019127 is decidedly less than our significance level of 0.05, we are statistically compelled to reject the **null hypothesis** (H_0). This rejection carries a critical statistical implication: we possess strong, sufficient evidence to conclude that the two population variances are, in fact, *not* equal. Practically speaking, this indicates that the inherent spread or measure of variability observed in the scores of population x is significantly different from the variability found in population y .

This conclusion holds significant consequences for any subsequent statistical analysis. For instance, if the next planned step was to perform a t-test to compare the mean scores of these two groups, the rejection of the equal **variance** assumption makes the standard independent samples t-test inappropriate. Instead, this finding mandates the use of a modified procedure, such as **Welch's t-test**, which is specifically designed to handle situations where the assumption of homoscedasticity has been violated.

Critical Assumptions and Operational Nuances of the F-Test

While the implementation of the **F-test** is computationally straightforward, the reliability of its results and the validity of any derived inferences are highly dependent on several fundamental assumptions concerning the data distribution. A deep understanding of these statistical nuances is essential for any professional data analysis.

The primary assumptions underpinning the F-test for the equality of variances include:

The two samples utilized in the test must have been drawn independently from their respective populations.

Each population from which the samples originate must be approximately **normally distributed**. The F-test is notoriously sensitive to violations of this normality assumption, particularly when dealing with smaller sample sizes. If the data cannot reasonably be assumed to be normal,

statisticians often recommend employing more robust, non-parametric alternatives, such as [Levene's test](#) or Bartlett's test, for assessing variance equality.

A crucial operational detail pertains to the structure of the custom Python function defined earlier. The standard convention in statistical practice dictates that the F-statistic should always be calculated by placing the sample with the larger variance (s_1^2) in the numerator and the sample with the smaller variance (s_2^2) in the denominator. This practice guarantees that the calculated F-statistic is always greater than or equal to 1 ($F \geq 1$).

If the input samples were defined arbitrarily, the function would ideally require additional internal logic to check which variance is larger and swap the arrays accordingly to adhere to this convention. Adhering to this structure is necessary to utilize the simplified one-tailed p-value calculation ($1 - \text{cdf}$) derived from the upper tail of the [F-distribution](#). Therefore, when calling the function, users must be diligent and ensure that the sample possessing the larger variance is passed as the first argument (x) to guarantee the F-statistic is correctly calculated and interpreted relative to the F-distribution based on the numerator and denominator [degrees of freedom](#).

The F-test remains a widely deployed tool used to answer specific investigative questions in fields ranging from quality control to medical research. Common applications include determining whether two distinct samples originate from populations possessing equivalent variability, or evaluating if a new intervention successfully reduces inherent variability compared to a standard baseline process. By correctly implementing the F-test in Python and meticulously adhering to its underlying assumptions, data scientists gain the confidence needed to determine whether observed differences in data dispersion are statistically significant, thereby making informed decisions regarding subsequent model selection and hypothesis validation.

Related: [How to Perform an F-Test in R](#)