

Hypothesis Testing

Task 5

Your last task is determining whether the average number of minutes watched in the US is similar to that in India.

Understanding the differences in usage patterns can help in product localization. The platform might need to tailor its content, features, or user interface to better fit the preferences or needs of users in different regions.

You'll focus only on free-plan students in 2022. Use the Excel sheet Task 5 to perform your calculations.

Your null hypotheses should (respectively) include the following:

- The engagement (minutes watched) in the US is higher than or equal to that in India ($\mu_1 \geq \mu_2$). We test only free-plan students.
- The engagement (minutes watched) in the US is lower than that in India ($\mu_1 < \mu_2$). We test only free-plan students.

Additionally, perform a two-sample t-test assuming unequal variances.

Optional: *Perform a two-sample f-test for variances to support the assumptions.*

What conclusion can you draw from this test? Is the engagement in the US higher than that in India?

Tip: Note that the degrees of freedom are calculated using the following formula for independent samples with unknown variances which are assumed to be unequal:

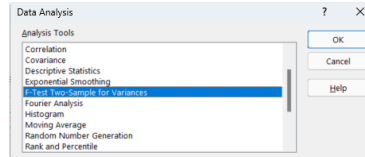
$$df = \frac{\left(\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y} \right)}{\frac{(s_x^2/n_x)^2}{n_x-1} + \frac{(s_y^2/n_y)^2}{n_y-1}}$$

Note: Assume that the degrees of freedom are equal to 11,001.

The t-test is equal to the following:

$$T = \frac{(\bar{x} - \bar{y}) - \mu_d}{\sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}}$$

First, you must perform a two-sample f-test for variances to prove that assumption of unequal variances between the samples (minutes watched by free-plan subscribers in the US and India):



Excel will perform the f-test for variances and provide the test statistic (f-value) and the p-value. The p-value would indicate the probability of obtaining the observed f-value if the null hypothesis (equal variances) were true.

As mentioned, compare the p-value to your chosen significance level (alpha) to determine if the variances are significantly different. If the p-value is less than or equal to your alpha level, you can reject the null hypothesis of equal variances.

The next step is to perform a t-test and compare it to the critical value from the t-distribution.

Consider the following steps for hypothesis testing where the variances are assumed unequal:

1. Specify the significance level:

$$\alpha = 1 - \text{Confidence Level} = 1 - 0.95 = 0.05$$

2. Calculate the t-statistic using the following formula:

$$T = \frac{(\bar{x} - \bar{y}) - \mu_d}{\sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}}$$

3. Look up the critical t-value using a t-distribution [table](https://365datascience.com/calculators/tables/std-table/) (<https://365datascience.com/calculators/tables/std-table/>) to correspond to your chosen significance level (commonly 0.05) and calculated degrees of freedom.

4. Note that the degrees of freedom are calculated using the following formula for independent samples with unknown variances which are assumed to be unequal:

$$df = \frac{\left(\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y} \right)}{\frac{(s_x^2/n_x)^2}{n_x-1} + \frac{(s_y^2/n_y)^2}{n_y-1}}$$

Note: Assume that the degrees of freedom are equal to 11,001.

5. Compare t-statistic to critical t-value. Interpret the magnitude and the sign of your t-statistic. The decision rule—based on the critical value approach—is as follows:

$$\text{If } T \leq -t_{df,0.05}, \text{ reject } H_0$$

$$\text{If } T > -t_{df,0.05}, \text{ fail to reject } H_0$$

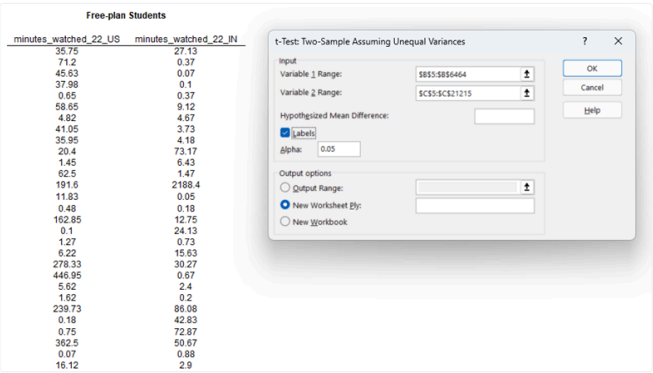
Assume that the critical t-value is equal to 1.65.

The decision rule—based on the p-value approach—is as follows:

$$p\text{-value} \leq \alpha, \text{ Reject } H_0$$

$$p\text{-value} > \alpha, \text{ Fail to reject } H_0$$

Alternatively, you can use the Data Analysis ToolPak in Excel to obtain the result directly:



Excel will perform the two-sample t-test assuming unequal variances and provide the results, including the t-statistic, degrees of freedom, and the p-value. The p-value would indicate the probability of obtaining the observed t-statistic if the null hypothesis (no difference between means) were true.

Compare the p-value to your chosen significance level (alpha) to determine if the difference between the means of the two samples is statistically significant. If the p-value is less than or equal to your alpha level, you can reject the null hypothesis of similar means.