

Table of Contents

1.Descriptive Statistics

- 1.1 Mean, Median, Standard Deviation
- 1.2 Skewness, Kurtosis
- 1.3 Conclusions

2. Confidence Interval

- 2.1 The process
- 2.2 The results

3. Hypothesis Testing I

- 3.1 Type I & Type II Errors

4. Hypothesis Testing II

The Case Study

In 2022, expectations were high for the growth of the 365 company and increased student engagement due to the introduction of new features on their website platform. These features included an XP system that allowed students to track progress, level up, and earn rewards by completing various learning objectives.

The platform also introduced in-app coins that could be exchanged for special awards, a leaderboard for students to compete for top positions in different divisions, earning weekly rewards and advancing up the ladder, and streaks to encourage consistent learning habits. Additionally, the company expanded its course library, offering a wider range of topics to equip students with diverse skills and attract a larger audience.

These enhancements were expected to improve the student experience, create an effective customer engagement strategy, and contribute to the company's success in the coming year. This Customer Engagement Analysis in Excel project requires you to analyze whether these new platform additions have increased student engagement.

1. Descriptive Statistics

1.1. Mean Median Standard Deviation

We will focus on low-engagement users (those who watched between 1 and 100 minutes in 2021). Low-engagement users often represent the most significant potential for growth. If 365 can find ways to increase its usage, it could significantly impact the overall use of the platform.

Paid-Plan Students		
	minutes_watched_21	minutes_watched_22
Mean	33.80	273.02
Median	26.33	40.28
Standard Deviation	28.21	854.58

- **Mean:** Among students who watched between 1 and 100 minutes in 2021, the average minutes watched by paid-plan students increased significantly from Q4 2021 to Q4 2022, from approximately 33.80 minutes to about 273.02 minutes. This suggests a substantial increase in engagement among this group of initially low-engagement-paid-plan students.
- **Median:** The median minutes these low-engagement-paid-plan students watched increased from Q4 2021 to Q4 2022, from 26.33 minutes to 40.28 minutes. While this increase is not as dramatic as the increase in the mean, it indicates that the typical student in this group (i.e., the student in the middle of the distribution) also increased their engagement. This suggests that the increase in engagement was more widespread among paid-plan students and not solely driven by a few outliers.
- **Standard Deviation:** The standard deviation for these low-engagement-paid-plan students increased substantially from 28.21 minutes in Q4 2021 to 854.58 minutes in Q4 2022. This indicates a much larger variability in the minutes watched by these students in Q4 2022 compared to Q4 2021. This could be due to a broader range of engagement levels among the students in Q4 2022, with some students watching very little content and others watching a lot of content.

These results suggest that paid-plan students who were initially low-engagement in 2021 significantly increased their engagement in 2022. But the increased standard deviation indicates a broader range of engagement levels among these students in 2022. Understanding the reasons behind this variability could provide valuable insights for further boosting engagement. For instance, the factors that motivated the students who significantly increased their engagement might be leveraged to encourage increased engagement among other students.

	Free-plan Students	
	minutes_watched_21	minutes_watched_22
Mean	25.39	117.64
Median	14.17	11.83
Standard Deviation	26.23	468.93

- **Mean:** Among students who watched between 1 and 100 minutes in 2021, the average minutes watched by free-plan students increased from about 25.39 minutes in Q4 2021 to about 117.64 minutes in Q4 2022. This suggests that overall engagement among these initially low-engagement-free-plan students increased during this period. But the extent of this increase is less than what was observed for similar low-engagement-paid-plan students, suggesting that while these free-plan students are watching more content, they're still not as engaged as the equivalent group of paid-plan students.
- **Median:** Interestingly, the median minutes watched by these low-engagement-free-plan students decreased from Q4 2021 to Q4 2022, from 14.17 minutes to 11.83 minutes. This indicates that engagement decreased for the typical student in this group (i.e., the student in the middle of the distribution). The increase in the mean might be driven by a small number of free-plan students who significantly increased their engagement in Q4 2022, while the majority did not increase their engagement or even reduced it.
- **Standard Deviation:** The standard deviation for the low-engagement-free-plan students increased from 26.23 minutes in Q4 2021 to 468.93 minutes in Q4 2022. This indicates a more significant variability in the minutes watched by these students in Q4 2022 compared to Q4 2021. The behavior of these students then became more diverse in Q4 2022, with some watching a lot of content and others watching very little.

These results suggest a complex picture for the initially low-engagement-free-plan students. While the mean minutes watched increased—signifying an increase in overall engagement—the median minutes

watched decreased, indicating that the typical student in this group did not increase their engagement. This discrepancy and the increased standard deviation suggest that a small number of students within this group might significantly increase their engagement while the majority did not. This might imply the need for targeted strategies to boost engagement among the broader population of initially low-engagement-free-plan students.

1.2. Skewness, Kurtosis

Skewness is a fundamental measure of probability distribution asymmetry in a dataset. It reveals whether the observations are concentrated more on one side of the distribution. This metric helps us understand how the data deviates from a normal distribution and provides insights into its underlying structure. A positive skewness value (higher than 0) indicates a right-skewed distribution, while a negative skewness value (lower than 0) points to a left-skewed distribution. A symmetrical distribution has a skewness value of 0, indicating a balanced data spread around the mean.

Kurtosis measures the degree of tailedness—the weight of the tails relative to the rest of the distribution. In other words, it shows how much of the data is in the tails compared to the center. Located farthest from the center, the tails represent the regions where data points are more dispersed — suggesting the presence of more extreme values. If a distribution is heavy-tailed — i.e., more data in the tails — it exhibits high kurtosis. Meanwhile, a low kurtosis occurs when the data is more evenly distributed between the tails and the center or the distribution is light-tailed. Kurtosis values greater than 0 indicate that the data has heavier tails and a sharper peak than the normal distribution (leptokurtic). A leptokurtic distribution has a high positive kurtosis, suggesting that it's very peaked and has a relatively large number of outliers. This type has a higher frequency of extreme values or outliers. Distributions with kurtosis less than 0 are called platykurtic. These distributions have lighter tails and a flatter peak compared to a normal distribution. This indicates that the data has fewer extreme values or outliers (platy kurtosis). If kurtosis equal to 0 are called mesokurtic. The normal distribution is an example of a mesokurtic distribution, where the tails and peak are neither particularly heavy nor light (zero kurtosis).

Paid-Plan Students			Free-Plan Students		
Skewness	0.63	7.07	Skewness	1.17	15.06
Kurtosis	-0.85	58.48	Kurtosis	0.36	315.76

1.3.Conclusions

- On average, low-engagement-paid students initially increased their watching time more significantly than the free-plan students from Q4 2021 to Q4 2022. This could suggest that paid-plan students find more value in the platform, possibly due to premium features or content that are available to them.

In contrast, the median watch time decreased for free-plan students, suggesting that the typical free-plan student in this group did not increase their engagement. This discrepancy might indicate that the strategies or features designed to increase engagement are more effective for paid-plan students. It could also suggest that the monetary investment leads to increased usage due to a desire to get their money's worth.

Based on the findings, the platform is more successful in increasing engagement among students who make a monetary investment (i.e., paid-plan students). But the increased variability, especially among paid-plan students, indicates that there are likely differences in how individual students are responding to the platform's offerings.

Therefore, personalized approaches might be beneficial in boosting engagement, and further analysis could help understand the factors that drive increased engagement among paid- and free-plan students.

- The skewness for **free-plan students** increased from 1.17 in Q4 2021 to 15.06 in Q4 2022, indicating positive skewness.

The mean is larger than the median in a right-skewed distribution because the distribution tail pulls the mean to the right. This observation is confirmed by the mean and median values in the two years. An increasing skewness suggests that more students watch significantly more content than most over time, pulling the mean upwards. In both cases, the mean is higher than the median ($33.80 > 26.33$ in 2021 and $273.02 > 40.28$ in 2022)

Overall, the increasing skewness and kurtosis for both groups from Q4 2021 to Q4 2022 suggest a growing number of students watching significantly more content than the majority. This is especially true for free-plan students with a higher skewness and kurtosis in Q4 2022 than paid-plan students.

2. Confidence Intervals

Confidence intervals are a statistical tool used to estimate the range within which a population parameter is likely to fall, providing an indication of the precision and reliability of the estimate. They are constructed based on sample data and offer a way to measure the degree of uncertainty or certainty in the estimation process.

The dataset we have to work with information about two pairs of students:

- Students who haven't had a paid-plan subscription (engaged in Q4 2021& engaged in Q4 2022)
- Students who have been paid-plan subscribers (engaged in Q4 2021& engaged in Q4 2022)

For each of the four groups, we'll determine the minute interval within which we can be 95% confident that a randomly selected individual will be situated. What conclusions can be drawn about the students' engagement in Q4 2021 and Q4 2022?

2.1. The Process

With a confidence level equal to 95%, we estimate the *significance level*

$$\alpha = 1 - \text{Confidence level}$$

Then, we use the sample's mean (\bar{x}), standard deviation (s) and size (n) to calculate the standard error (SE).

$$s_{\bar{x}} = \frac{s}{\sqrt{n}}$$

Since the sample size is larger than 30, we can assume normality. The z-score for a 95% confidence interval is 1.96. With the following formula we calculate the margin of error (ME).

$$ME = z_{\alpha/2} \cdot \frac{s}{\sqrt{n}}$$

Finally, the confidence interval is equal:

$$\bar{X} \pm ME$$

2.2. The results

Paid-Plan Students						
minutes_watched_21	minutes_watched_22					
2973.67	4110.17					
2939.48	4099.42					
2860.78	4085.2					
2853.73	4064.35					
2830.2	4024.33					
2809.67	3948.85					
2803.17	3909.85					
2797.55	3908.57					
2782.08	3879.82					
2741.9	3866.8					
2703.03	3828.88					
2699.63	3776.67					

- **Paid – Plan Students**

For paid-plan students, there's an increase in engagement from Q4 2021 to Q4 2022. This suggests that we can be 95% confident that the true average minutes watched by all paid-plan students in the population increased from Q4 2021 to Q4 2022.

- **Free – Plan Students**

Among free-plan students, there's a decrease in engagement from Q4 2021 to Q4 2022. We then can be 95% confident that the true average minutes watched by all free-plan students in the population decreased from Q4 2021 to Q4 2022.

- **Comparison between Paid – Plan & Free – Plan**

Students with a paid-plan subscription watch significantly more content than those without. In Q4 2022, the confidence interval for the average minutes watched by free-plan students was between 61.71 and 70.59 minutes, while for paid-plan students, it was between 351.99 and 384.72 minutes. This means we can be 95% confident that paid-plan students watched substantially more minutes than free-plan students in Q4 2022. This observation supports the expectation that paid-plan students, having invested in the platform, are more engaged than free-plan users.

Note that these interpretations are based on the confidence intervals. Further analysis is required to establish actual cause-effect relationships behind these engagement changes.

The observation that paid-plan subscribers watch more does not necessarily mean that having a paid-plan subscription directly encourages them to watch more. Higher engagement among paid-plan students may result from additional features or content available to them or because more engaged students are more likely to opt for a paid-plan subscription.

Similarly, the decrease in engagement among free-plan students could be attributed to various factors, such as changes in the platform, competition from other platforms, or changes in the user base.

3. Hypothesis Testing I

We want to reach a data-driven customer engagement decision on whether the platform's new features contribute to the increase of minutes watched on the platform for both free-plan and paying students — i.e., the rise in student engagement in their study process. To do that, we'll use hypothesis testing on both groups (free-plan and paying) for 2021 and 2022.

The null hypotheses (H_0) will include the following:

- The engagement (minutes watched) in Q4 2021 is higher than or equal to the one in Q4 2022 ($\mu_1 \geq \mu_2$). We test free-plan and paying students separately.

First, we conduct research on whether the variances between our samples are different, using a “Two-Sample F-Test for Variances”.

Two-Sample F-Test for Variances

Paid Students

	<i>minutes_watched_21</i>	<i>minutes_watched_22</i>
Mean	332.502508	368.3547139
Variance	236063.3116	355699.1148
Observations	3433	5104
df	3432	5103
F	0.663660104	
P(F<=f) one-tail	0	
F Critical one-tail	0.949796198	

Free Students

	<i>minutes_watched_21</i>	<i>minutes_watched_22</i>
Mean	133.9333129	69.14765544
Variance	134881.7038	65343.34428
Observations	32171	120658
df	32170	120657
F	2.06419958	
P(F<=f) one-tail	0	
F Critical one-tail	1.014667161	

The p-value indicates the probability of obtaining the observed f-value if the null hypothesis (equal variances) were true. The sample variances are not identical since the p-value in both cases is 0. Next, we use a left-tailed t-test assuming unequal variances for paying and free-plan students.

Now again using the Data Analysis Tool pack from Excel, we execute “Two-Sample t-Test Assuming Unequal Variances”.

Decision Rule: If $p - value \leq 0.05$ Reject H_0

Paid-Plan Students			
minutes_watched_21	minutes_watched_22	minutes_watched_21	minutes_watched_22
2973.67	4110.17	Mean	332.50
2939.48	4099.42	Standard Deviation	485.86
2860.78	4085.2	Sample Size	3433
2853.73	4064.35		
2830.2	4024.33		
2809.67	3948.85		
2803.17	3909.85		
2797.55	3908.57		
2782.08	3879.82		
2741.9	3866.8		
2703.03	3828.88		
2699.63	3776.67		
2686.85	3774.22		
2651.47	3754.5		
2649.98	3726.42		
2631.4	3699.57		
2574.6	3614.72		
2573.95	3607.25		
2571.68	3589.12		

Paid Students		minutes_watched_21	minutes_watched_22
Mean		332.502508	368.3547139
Variance		236063.3116	355699.1148
Observations		3433	5104
Hypothesized Mean Differen		0	
df		8229	
t Stat		-3.046942872	
P(T<=t) one-tail		0.001159572	
t Critical one-tail		1.645038819	
P(T<=t) two-tail		0.002319144	
t Critical two-tail		1.960252308	

Conclusion: Reject the null hypothesis because the p-value is lower than the specified significance level α (0.05).

Summary: With a t-statistic of -3.05 (less than the critical value of -1.645), you would reject the null hypothesis, indicating that the mean minutes watched by students in Q4 2021 is significantly smaller than the mean minutes watched by students in Q4 2022. This contradicts the null hypothesis, leading to its rejection. However, rejecting the null hypothesis does not confirm the alternative hypothesis; it simply suggests that the data provide sufficient evidence against the null hypothesis.

Free-plan Students			
minutes_watched_21	minutes_watched_22		
4716.68	6338.07		
4670.7	6280.12		
4622.35	6250.32		
4617.75	6208.8		
4599.53	6204.55		
4581.45	6099.35		
4568.73	6091.17		
4564.67	6073.37		
4481.85	6072.77		
4464.67	6071.8		
4439.47	6068.17		
4388.53	6043.12		
4385.1	5998.4		
4131	5949.77		
4123.17	5918.68		
4119.53	5875.6		
4115.7	5852.32		
4114.38	5849.75		
4112.93	5821.23		

	minutes_watched_21	minutes_watched_22
Mean	133.93	69.15
Standard Deviation	367.26	255.62
Sample Size	32171	120658

Free Students		
	minutes_watched_21	minutes_watched_22
Mean	133.9333129	69.14765544
Variance	134881.7038	65343.34428
Observations	32171	120658
Hypothesized Mean Difference	0	
df	40836	
t Stat	29.77523819	
P(T<=t) one-tail	4.7441E-193	
t Critical one-tail	1.644890942	
P(T<=t) two-tail	9.4881E-193	
t Critical two-tail	1.960022079	

Conclusion: For free-plan students: With a t-statistic of 29.78 (greater than the critical value of -1.645), you would fail to reject the null hypothesis. This means there's not enough evidence to conclude that μ_1 is smaller than μ_2 . So, the data supports the null hypothesis that μ_1 is larger than or equal to μ_2 .

These results align with previous findings from the confidence intervals and further underscore the difference in engagement patterns between paid- and free-plan students.

3.1 Type I & Type II Errors

A *Type I error* (false positive) occurs when you reject the null hypothesis even though it is true. In this context, it would mean concluding that engagement in 2022 is higher when it is not. The probability of making this error is represented by the significance level, α . Since the researcher selects the significance level for the hypothesis test, the responsibility for this error lies with the researcher.

The significance level is directly related to the confidence level, which represents our degree of certainty in the estimated results. It is equal to $(1 - \alpha)$. For instance, a 5% significance level in a hypothesis test implies a 5% probability of rejecting a true null hypothesis, corresponding to a 95% confidence level.

A *Type II error* (false negative) occurs when you fail to reject the null hypothesis when it is actually false. In this scenario, it would mean incorrectly concluding that engagement in 2022 is not higher when it actually is.

Which type of this error would cost more to the company?

The cost to the company for each type of error depends on the consequences of incorrect conclusions. If the company incorrectly concludes that engagement has increased (Type I error), it may over-invest in certain features or become complacent about the need for further improvements. Conversely, if the company incorrectly concludes that engagement has not increased (Type II error), it might fail to recognize successful features or miss opportunities to identify areas needing improvement.

4. Hypotheses test II

Our final task is to determine whether the average number of minutes watched in the US is similar to that in India.

Understanding the differences in usage patterns can help us with 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. We will focus only on free-plan students in 2022.

Our null hypotheses 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 will test only free-plan students.
- The engagement (minutes watched) in the US is lower than that in India ($\mu_1 < \mu_2$). We will test only free-plan students.

Decision Rule: If p – value ≤ 0.05 , reject H_0

Free-plan Students			
minutes watched 22 US	minutes watched 22 IN		
35.75	27.13	Mean	73.07
71.2	0.37	Standard Deviation	308.56
45.63	0.07	Sample Size	6459
37.98	0.1		
0.65	0.37		
58.65	9.12		
4.82	4.67		
41.05	3.73		
35.95	4.18		
20.4	73.17		
1.45	6.43		
62.5	1.47		
191.6	2188.4		
11.83	0.05		
0.48	0.18		
162.85	12.75		
0.1	24.13		
1.27	0.73		
6.22	15.63		
278.33	30.27		
446.95	0.67		

Free-plan Students	
minutes watched 22 US	minutes watched 22 IN
Mean	73.07053569
Variance	95208.64187
Observations	6459
Hypothesized Mean Difference	0
df	11001
t Stat	-1.210387573
P(T<=t) one-tail	0.113078106
t Critical one-tail	1.644992151
P(T<=t) two-tail	0.226156213
t Critical two-tail	1.960179649

Conclusion: Fail to reject the null hypothesis because the p-value is higher than the specified significance level α (0.05).

If we had rejected the hypothesis that US students watch more or an equal amount of content as Indian students, it would suggest that US students watch less content on average than students in India.

This could have several implications:

1. **Market Differences:** This might indicate that the platform is more engaging or relevant to students in India than to US students. Understanding why this is the case could be valuable. Are specific features, content, or aspects of the platform particularly appealing to Indian students? These questions need further exploration but are beyond the scope of this analysis.
2. **Growth Opportunities:** Lower engagement among US students could represent a growth opportunity. The 365 company might explore ways to increase engagement among US students, such as through targeted marketing efforts, adding more relevant content, or other strategies.
3. **Resource Allocation:** This information could be useful in deciding where to allocate resources. For instance, if Indian students are more engaged, it might make sense to invest more in content and features targeted toward this audience.