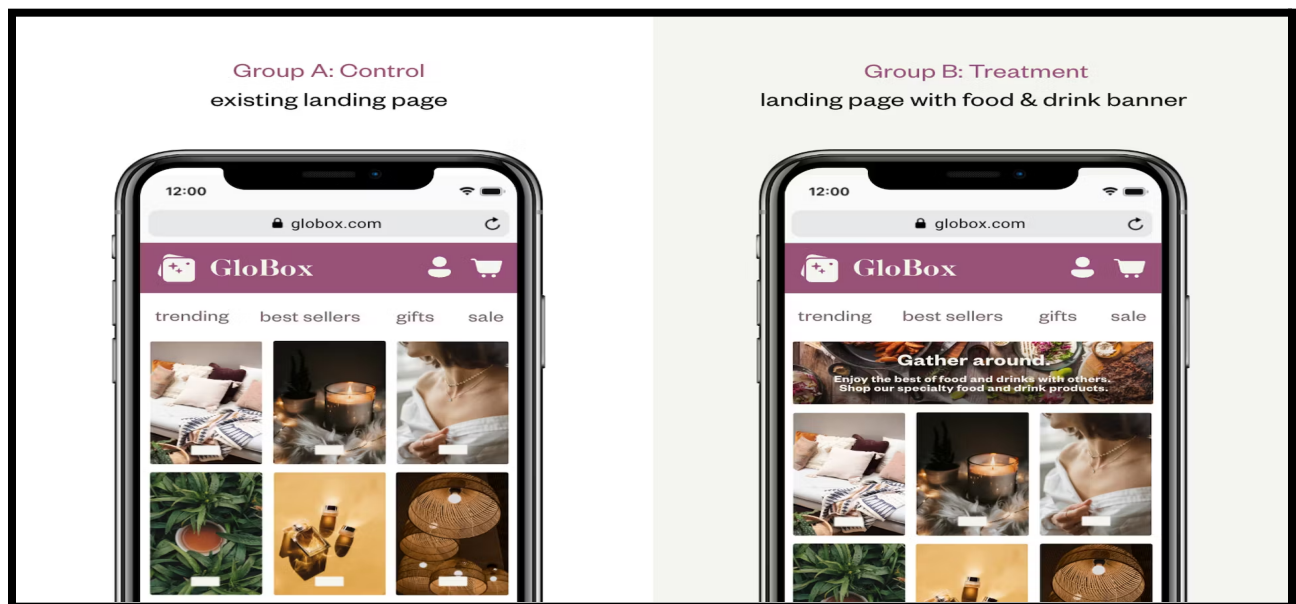# GloBox

## An e-commerce company

## A/B Testing Analysis of new webpage

### Company Overview

GloBox is an online marketplace that specializes in sourcing unique and high-quality products from around the world.GloBox is primarily known amongst its customer base for boutique fashion items and high-end decor products. However, their food and drink offerings have grown tremendously in the last few months, and the company wants to bring awareness to this product category to increase revenue.

The Growth team decides to run an A/B test that highlights key products in the food and drink category as a banner at the top of the website. The control group does not see the banner, and the test group sees it as shown below:
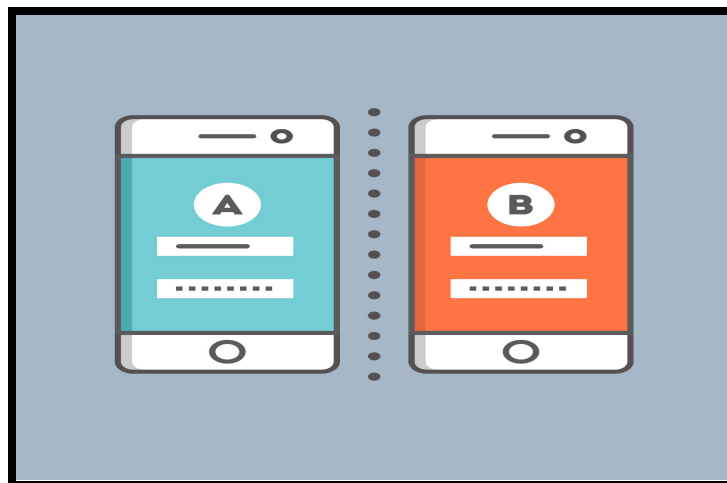
**The setup of the A/B test is as follows:**

1. The experiment is only being run on the mobile website.

2. A user visits the GloBox main page and is randomly assigned to either the control or test group. This is the join date for the user.

3. The page loads the banner if the user is assigned to the test group, and does not load the banner if the user is assigned to the control group.

4. The user subsequently may or may not purchase products from the website. It could be on the same day they join the experiment, or days later. If they do make one or more purchases, this is considered a "conversion".

**Task**

To analyze the results of the A/B test and provide a recommendation to our stakeholders about whether GloBox should launch the experience to all users.
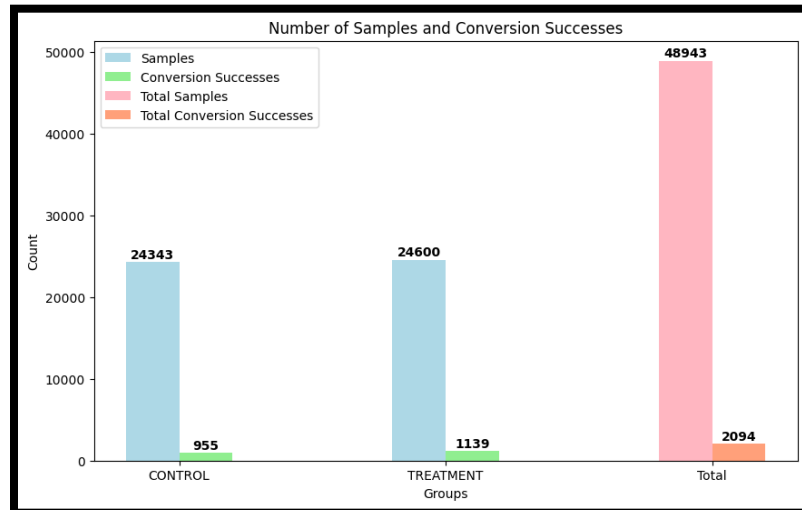
## Summary

I recommend that we do launch the new homepage because we did observe strong evidence that there was an increase in conversion rate during the stipulated period when the test was run. The new homepage with foods and drinks offerings would substantially increase the revenue.

**Analysis:**

★ Total no of sample users in each group and conversion successes were:

To demonstrate the visual understanding of the above conclusion we used python to get the graph plotted with the following code:



★ There were **767** Users in the control group in Canada.


The

★ Conversion rate for all users were **4.28%**

★ There were total **41412** users as of February 1st 2023 in the A/B test



★ The average amount spent per user for the control and treatment groups were:

★ The total spent on each group and percentage share:

To demonstrate the visual understanding of the above conclusion we used python to get the graph plotted with the following code:



★ The 95% confidence interval for the average amount spent per user in the control (used the t distribution)

★ The 95% confidence interval for the average amount spent per user in the treatment. (Used the t distribution)



To demonstrate the visual understanding of the above conclusion we used python to get the graph plotted with the following code:

| Groups | Sample | Mean | STDV | Spent | No of Conversion Successes | Conversions | Conversion rate in % | Difference in Conversion rates: (treatment - control) | 95% confidence interval for the average amount spent per user | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | Lower Bound | Upper Bound |
| CONTROL | 24343 | 3.375 | 25.93639056 | 82146 | 955 | 0.03923099043 | 3.92% | 0.00706982258 | 3.04869 | 3.70035 |
| TREATMENT | 24600 | 3.391 | 25.4141096 | 83415 | 1139 | 0.04630081301 | 4.63% | | 3.07327 | 3.70846 |
| Total | 48943 | 3.383 | 25.67494579 | 165561 | 2094 | 0.04278446356 | 4.28% | | | |

★ To Conduct a hypothesis test to see whether there is a difference in the average amount spent per user between the two groups. (Used the t distribution and a 5% significance level. Assuming unequal variance.)

We used this code to download the required columns as csv file to further process the query in google sheets.



The resulting p-value and conclusion derived from the sheet calculations:



| **Q-4:** Conduct a hypothesis test to see whether there is a difference in the average amount spent per user between the two groups. What are the resulting p-value and conclusion? Use the t distribution and a 5% significance level. Assume unequal variance. | | | | |
|---|---|---|---|---|
| H0: μ1 = μ2 | Control Mean (AVERAGE(C2:C24345)) | 3.374518468 | Treatment Mean (AVERAGE(F2:F)) | 3.390866946 |
| Ha: μ1 ≠ μ2 | Control Stdev (STDEV(C2:C24345)) | 25.936390557 | Treatment Stdev (STDEV(F2:F)) | 25.414109599 |
| | n1-sample control | 24343 | n2-sample treatment | 24600 |
| | T Test | 0.943384262 | | |

p = 0.944, statistically _insignificant_. We _fail to reject_ the null hypothesis that there is no difference in the mean amount spent per user between the control and treatment. We are using a two-sided t-test for a difference in means. Assuming unequal variance, we use the unpooled standard error.

To demonstrate the visual understanding of the above conclusion we used python to get the graph plotted with the following code:

Based on the given options and assuming a 5% significance level, p = 0.944, statistically insignificant. We reject the null hypothesis that there is no difference in the mean amount spent per user between the control and treatment.

Since the p-value calculated from your t-test is 0.9434 and it is greater than the significance level of 0.05, we cannot reject the null hypothesis that there is no difference in the mean amount spent per user between the two groups. However, if the significance level were higher or if the p-value were lower, we would have sufficient evidence to reject the null hypothesis and conclude that there is a statistically significant difference in the mean amount spent per user between the two groups.

★ The 95% confidence interval for the difference in the average amount spent per user between the treatment and the control (treatment-control) using the t distribution and assumed unequal variance.

| Q-5: What is the 95% confidence interval for the difference in the average amount spent per user between the treatment and the control (treatment-control)? Use the t distribution and assume unequal variance. | |
|---|---|
| standard error of the difference between the two sample means: SQRT((K3^2/K4)+(I3^2/I4)) | 0.232140559 |
| margin of error: T.INV.2T(0.025,K4+I4-2)*I9 | 0.520336493 |
| Upper bound: (mean_b - mean_a) + 1.96*se | 0.471343973 |
| Lower bound: (mean_b - mean_a) - 1.96*se | -0.438647017 |
| (-0.439, 0.471) We are using a two-sample t-interval for a difference in means. Assuming unequal variance, we use the unpooled standard error. | |

(-0.439, 0.471) are the confidence intervals where we are using a two-sample t-interval for a difference in means. Assuming unequal variance, we use the unpooled standard error.

To demonstrate the visual understanding of the above conclusion we used python to get the graph plotted with the following code:



★ The user conversion rate for the control and treatment groups



| Groups | Sample | No of Conversion Successes | Conversions | Conversion rate in % | Mean | STDV | Difference in Conversion rates: (treatment - control) |
|---|---|---|---|---|---|---|---|
| CONTROL | 24343 | 955 | 0.03923099043 | 3.92% | 3.374518468 | 25.93639056 | 0.00706982258 |
| TREATMENT | 24600 | 1139 | 0.04630081301 | 4.63% | 3.390866946 | 25.4141096 | |

To demonstrate the visual understanding of the above conclusion we used python to get the graph plotted with the following code:

★ The 95% confidence interval for the conversion rate of users in the control using the normal distribution.

| Q-7: What is the 95% confidence interval for the conversion rate of users in the control? Use the normal distribution. | | |
|---|---|---|
| Conversion rate of control group - "A" | 0.039230990 | |
| The standard error (SE) for the sample proportion in the control group would be: = sqrt(p*(1-p)/n) | 0.001244334 | |
| Sample size: n = 24343 | | |
| Sample proportion (control group): p = 0.0392 | | |
| Critical value for 95% confidence level: z* = 1.96 | | |
| CI = p ± z* * SE (Upper bound) | 0.04167 | |
| CI = p ± z* * SE (Lower bound) | 0.036792095 | |
| Therefore, we can be 95% confident that the true conversion rate of users in the control group is between 3.68% and 4.16%. We are using a one-sample z-interval for proportions. | | |

Therefore, we can be 95% confident that the true conversion rate of users in the control group is between 3.68% and 4.17%. We are using a one-sample z-interval for proportions.

To demonstrate the visual understanding of the above conclusion we used python to get the graph plotted with the following code:
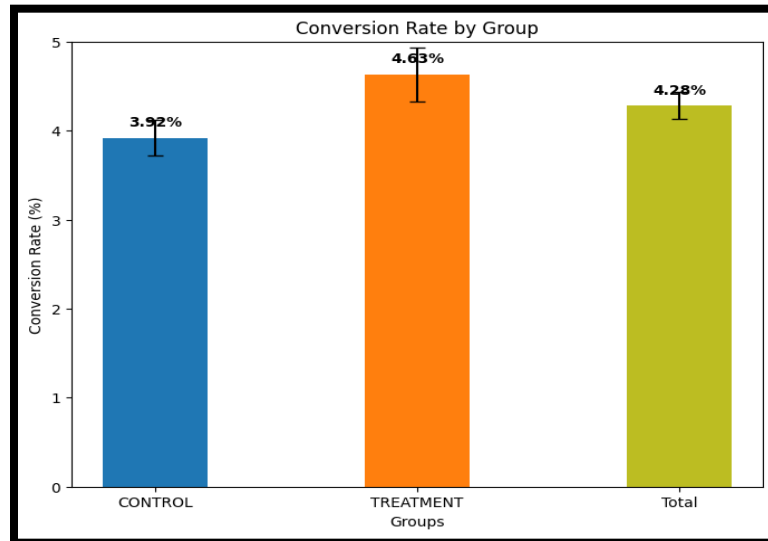
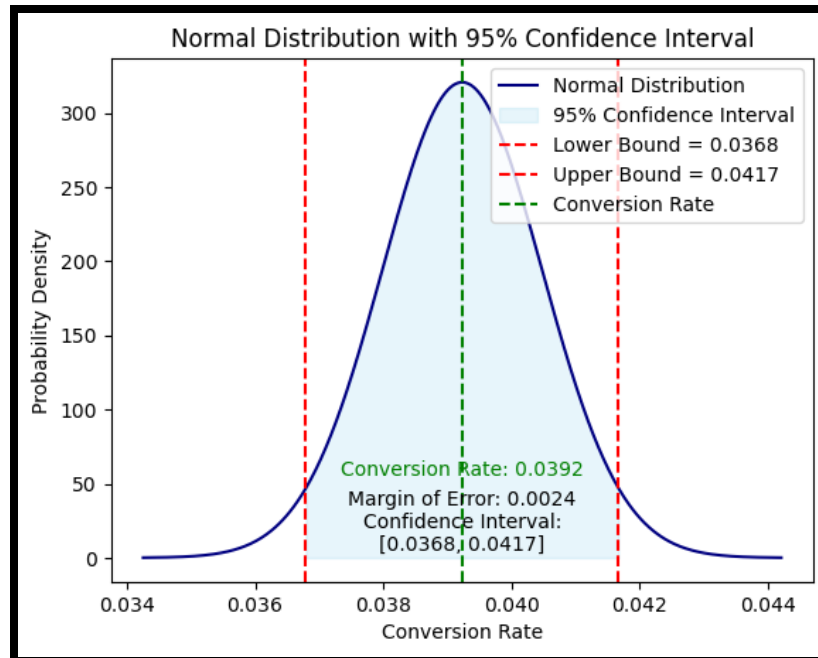Normal Distribution with 95% Confidence Interval

- ★ The 95% confidence interval for the conversion rate of users in the treatment using the normal distribution.

| Q-8: What is the 95% confidence interval for the conversion rate of users in the treatment? Use the normal distribution. | | |
|---|---|---|
| Conversion rate of Treatment group - "B" | 0.046300813 | |
| The standard error (SE) for the sample proportion in the control group would be: = sqrt(p*(1-p)/n) | 0.001339777 | |
| Sample size: n = 24600 | | |
| Sample proportion (control group): p = 0.0392 | | |
| Critical value for 95% confidence level: z* = 1.96 | | |
| CI = p ± z* * SE (Upper bound) | 0.048926776 | |
| CI = p ± z* * SE (Lower bound) | 0.043674850 | |
| Therefore, we can be 95% confident that the true conversion rate of users in the treatment group is between 4.37% and 4.89%. We are using a one-sample z-interval for proportions. | | |

Therefore, we can be 95% confident that the true conversion rate of users in the treatment group is between 4.37% and 4.89%. We are using a one-sample z-interval for proportions.

To demonstrate the visual understanding of the above conclusion we used python to get the graph plotted with the following code:
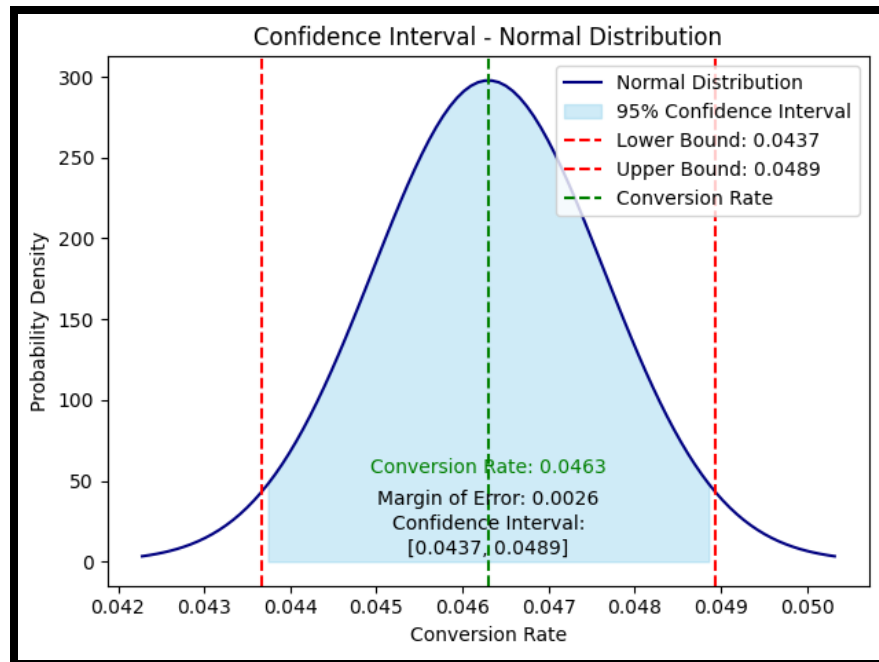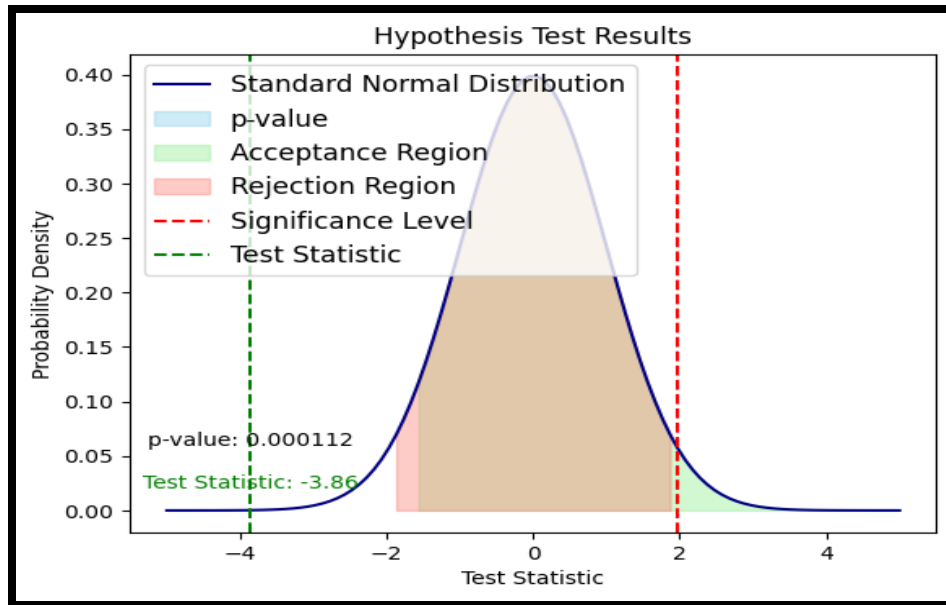
Confidence Interval - Normal Distribution

Conversion Rate: 0.0463
Margin of Error: 0.0026
Confidence Interval:
[0.0437, 0.0489]

★ To Conduct a hypothesis test to see whether there is a difference in the conversion rate between the two groups using the normal distribution and a 5% significance level. Taken into consideration the pooled proportion for the standard error.

**Q-9:** Conduct a hypothesis test to see whether there is a difference in the conversion rate between the two groups. What are the resulting p-value and conclusion? Use the normal distribution and a 5% significance level. Use the pooled proportion for the standard error.

| | CONTROL | TREATMENT |
|---|---|---|
| H0: p1 - p2 = 0 | | |
| Ha: p1 - p2 ≠ 0 | X1 | X2 |
| The number of successes represents the number of users who converted in each group. In the context of a hypothesis test, a success is defined as the event of interest, such as a user making a purchase, signing up for a service, or clicking on a button. | 955 | 1139 |
| The Pooled proportion: p^ = x1+x2/n1+n2 // OR | 0.042784464 | |
| P^ = (p1*n1 + p2*n2) / (n1 + n2) | 0.0427845 | |
| where p1 is the conversion rate of the control group, and p2 is the conversion rate of the treatment group. | | |
| The standard error: | 0.001829526 | |
| SE = sqrt(P^ * (1 - P^) * ((1/n1) + (1/n2))) | | |
| The test statistic is t = (p1 - p2) / SE | -3.864291770 | |
| The p-value can be calculated using a two-tailed t-distribution with degrees of freedom equal to n1 + n2 - 2: | | |
| df = n1 + n2 - 2 | 48941 | |
| Using this degrees of freedom, the p-value can be calculated as: | | |
| p-value = 2 * T.DIST(t, df, 1) | 0.000111556 | |

Since the p-value is less than 0.05, i.e. p = 0.0001, we can reject the null hypothesis and conclude that there is evidence of a statistically significant difference in conversion rates between the control and treatment groups. We are using a two-sample two-sided z-interval for a difference in proportions. Assuming equal proportions, we use the pooled standard error.

Since the p-value is less than 0.05, i.e. p = 0.0001, we can reject the null hypothesis and conclude that there is evidence of a statistically significant difference in conversion rates between the control and treatment groups. We are using a two-sample two-sided z-interval for a difference in proportions. Assuming equal proportions, we use the pooled standard error.

To demonstrate the visual understanding of the above conclusion we used python to get the graph plotted with the following code:

Hypothesis Test Results

★ The 95% confidence interval for the difference in the conversion rate between the treatment and control (treatment-control) using the normal distribution and unpooled proportions for the standard error.
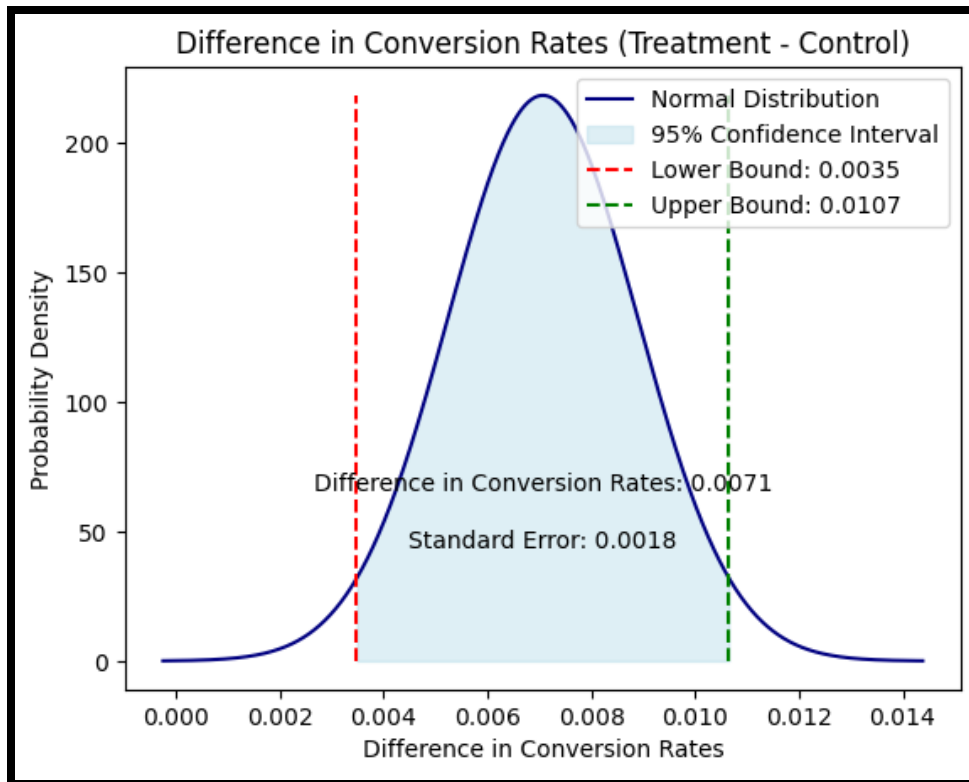
**Q-10:** What is the 95% confidence interval for the difference in the conversion rate between the treatment and control (treatment-control)? Use the normal distribution and unpooled proportions for the standard error.

| | |
|---|---|
| Difference in Conversion rates: (treatment - control) | 0.007069823 |
| Using unpooled proportions for the standard error, we can calculate the standard error as: | |
| SE = sqrt(p1*(1-p1)/n1 + p2*(1-p2)/n2) | 0.001828488 |
| To calculate the 95% confidence interval, we can use the formula: | |
| CI = (p1 - p2) ± z*SE | |
| Upper Bound | 0.010653660 |
| Lower Bound | 0.003485985 |

Therefore, we can say with 95% confidence that the true difference in conversion rates between the treatment and control lies between 0.0035 and 0.0107. Since the interval does not contain zero, we can conclude that there is a statistically significant difference in the conversion rates between the treatment and control groups.

Therefore, we can say with 95% confidence that the true difference in conversion rates between the treatment and control lies between 0.0035 and 0.0107. Since the interval does not contain zero, we can conclude that there is a statistically significant difference in the conversion rates between the treatment and control groups.

To demonstrate the visual understanding of the above conclusion we used python to get the graph plotted with the following code:

### Recommendation

Based on the results above, it does make sense to launch the treatment because we did observe an increase in conversion rate per user. I recommend that we do launch it. So the significant surge on the Conversion rates do give confidence that our revenue would also increase significantly with the introduction of foods and drinks on the new homepage. On the other hand I also observed that the difference of p-value with the significance level is very small and understand if possibility of increasing the sample size is there than we can relaunch the test to get better results.

### Appendices

➢ **Analysis in Python - Code file attached as zip folder.**
➢ **Analysis in SQL - Code file attached as zip folder.**
➢ **⊞ DA201 Mastery Projects- A/B Testing Analysis for an e-commerce company called GloBox.csv**
➢ **▭ DA201 Mastery Projects- A/B Testing Presentation for an e-commerce company called GloBox**