



MASTERY PROJECT 1

Fereshteh Ranjbar



MAY 12, 2023
MASTER SCHOOL

Table of Contents

<i>Table of Figures</i>	2
<i>Introduction</i>	3
<i>Calculations</i>	4
Conversion Rate	4
Group A	4
Control Group's Conversion Rate (A/Old version).....	4
Conversion Rate Calculation Results	4
Control Group's Conversion Rate 95% Confidence Interval (A/Old version).....	5
Confidence Interval Results.....	5
Group B	6
Treatment Group's Conversion Rate (B/New version)	6
Conversion Rate Calculation Results	6
Treatment Group's Conversion Rate 95% Confidence Interval (B/New version)	7
Confidence Interval Results.....	7
<i>Comparison Between Groups</i>	8
Pooled Probability of Conversion	8
SQL code & Results	8
Pooled Standard Error	9
SQL code & Results	9
<i>Hypothesis Test</i>	10
Difference Calculation	10
<i>Data Visualization</i>	11
<i>Conclusion & Recommendations</i>	13

Table of Figures

Figure 1 Conversion rate for group A.....	4
Figure 2 Results for group A.....	4
Figure 3 95% Confidence interval calculation for group A	5
Figure 4 Confidence interval results group A	5
Figure 5 Conversion rate for group A.....	6
Figure 6 Results for group B.....	6
Figure 7 Confidence interval for group B.	7
Figure 8 Confidence interval results for group B.....	7

Introduction

The purpose of this analysis is to evaluate and compare two variations of the app/website to determine the most viable course of action, whether it is launching the new version, refraining from launching, or re-evaluating the decision to change the front end.

This decision will be based on a thorough examination of conversion rates, which serve as a key metric for assessing user engagement and success. By analyzing and comparing the conversion rates between the two variations, we can derive valuable insights regarding user behavior and preferences, allowing us to make an informed decision regarding the future direction of the app/website.

Calculations

In this section, we will present the calculations, SQL code, and visualizations that were used to analyze the data. By employing these techniques, we aim to gain a comprehensive understanding of the underlying patterns and trends.

Conversion Rate

To begin the analysis, we extracted the relevant data from the database using SQL queries.

1. Conversion Rate Calculation: I computed the conversion rate by dividing the number of successful conversions by the total number of visitors or users. The formula used was:

$$\text{Conversion Rate} = \left(\frac{\text{Successful Conversions}}{\text{Total Visitors}} \right) * 100$$

$$\text{Conversion Rate} = \frac{\text{Number of users in "activity" table}}{\text{Number of users in the "group" table}}$$

Group A

This section includes all the necessary calculations for coming to a reasonable conclusion for group A.

Control Group's Conversion Rate
(A/Old version)

The following is the SQL code for calculating the conversion rate for the control group (A).

Figure 1 Conversion rate for group A.

```
1  with control as(
2  |   select count(distinct uid) as numcont
3  |   from groups as g
4  |   where g.group='A'
5  | ),
6  numactA as(
7  |   select count(distinct activity.uid)
8  |   from activity
9  |   left join groups
10 |   on activity.uid=groups.uid
11 |   where groups.group = 'A'
12 |
13 |   SELECT CAST((SELECT * FROM numactA) AS FLOAT) /
14 |   (SELECT * FROM control) AS controlconvrate;
```

Conversion Rate Calculation Results

The following is the result of the SQL query which shows the conversion rate (in percentage)

Query Results	
	1 ROWS
controlconvrate	FLOAT8
3.9230990428459926	

Figure 2 Results for group A.

Control Group's Conversion Rate 95% Confidence Interval (A/Old version)

```
confidence_level = 0.95 # 95% confidence level
Z_value = 1.96 # Z value for a 95% confidence interval
lower_bound_A = conversion_rate_A - Z_value * standard_error_A
upper_bound_A = conversion_rate_A + Z_value * standard_error_A
confidence_interval_A = (lower_bound_A, upper_bound_A)
```

The code snippet provided in Figure 2 demonstrates the calculation of the 95% confidence interval for the conversion rate of group A in a dynamic manner.

It is important to note that the Z value of 1.96, obtained from the Z-value table, is used in this calculation. Also, the equation box below is the theoretical steps behind calculating the interval.

```
1 ∵ WITH control AS(
2     SELECT COUNT(DISTINCT uid) AS numcont
3     FROM groups AS g
4     WHERE g.group='A'
5 ),
6 ∵ numactA AS(
7     SELECT COUNT(DISTINCT activity.uid)
8     FROM activity
9     LEFT JOIN groups ON activity.uid=groups.uid
10    WHERE groups.group = 'A'
11 ),
12 ∵ conversion_rate AS (
13     SELECT CAST((SELECT * FROM numactA) AS FLOAT) / (SELECT * FROM control) AS rate
14 ),
15 ∵ se AS (
16     SELECT SQRT(rate*(1-rate)/control.numcont) AS se
17     FROM conversion_rate, control
18 )
19 SELECT (rate - 1.96*se)*100 AS lower_bound, (rate + 1.96*se)*100 AS upper_bound
20 FROM conversion_rate, se;
```

Figure 3 95% Confidence interval calculation for group A.

Confidence Interval Results

The following is the result of the SQL query which shows the conversion rate 95% confidence interval.

Query Results	
1 ROWS	
lower_bound	upper_bound
FLOAT8 3.679209516027415	FLOAT8 4.1669885696645705

Figure 4 Confidence interval results group A.

Group B

This section includes all the necessary calculations for coming to a reasonable conclusion for group B.

Treatment Group's Conversion Rate (B/New version)

The following is the SQL code for calculating the conversion rate for the treatment group (B).

```
1  with treat as(
2      select count(distinct uid) as numcont
3      from groups as g
4      where g.group='B'
5  ),
6  numactB as(
7      select count(distinct activity.uid)
8      from activity
9      left join groups
10     on activity.uid=groups.uid
11     where groups.group = 'B'
12  )
13  SELECT CAST((SELECT * FROM numactB) AS FLOAT)*100 /
14  (SELECT * FROM treat) AS treatconvrate;
```

Figure 5 Conversion rate for group A

Conversion Rate Calculation Results

The following is the result of the SQL query which shows the conversion rate (in percentage)

Query Results	
1 ROWS	
controlconvrate	
FLOAT8	4.630081300813008

Figure 6 Results for group B.

Treatment Group's Conversion Rate 95% Confidence Interval (B/New version)

The code snippet provided in Figure 5 demonstrates the calculation of the 95% confidence interval for the conversion rate of group B in a dynamic manner.

```
1  WITH treat AS(
2    SELECT COUNT(DISTINCT uid) AS numcont
3    FROM groups AS g
4    WHERE g.group='B'
5  ),
6  numactB AS(
7    SELECT COUNT(DISTINCT activity.uid)
8    FROM activity
9    LEFT JOIN groups ON activity.uid=groups.uid
10   WHERE groups.group = 'B'
11  ),
12  conversion_rate AS (
13    SELECT CAST((SELECT * FROM numactB) AS FLOAT) / (SELECT * FROM treat) AS rate
14  ),
15  se AS (
16    SELECT SQRT(rate*(1-rate)/treat.numcont) AS se
17    FROM conversion_rate, treat
18  )
19  SELECT (rate - 1.96*se)*100 AS lower_bound, (rate + 1.96*se)*100 AS upper_bound
20  FROM conversion_rate, se;
```

Figure 7 Confidence interval for group B.

Confidence Interval Results

The following is the result of the SQL query which shows the conversion rate 95% confidence interval.

Query Results	
1 ROWS	
lower_bound	upper_bound
FLOAT8	FLOAT8
4.367485043175266	4.8926775584507505

Figure 8 Confidence interval results for group B.

Comparison Between Groups

This section encompasses the calculation of crucial values required for conducting a hypothesis test, which forms the basis for drawing a reasonable conclusion. By performing this analysis, we aim to assess the significance of any differences between the groups under investigation. The calculated values play a pivotal role in guiding our decision-making process.

Pooled Probability of Conversion

The pooled probability is the total probability of conversion across groups.

$$P_{pool} = \frac{\text{Successful conversion}_A + \text{Successful conversion}_B}{\text{Total}_A + \text{Total}_B}$$

SQL code & Results

The following is the SQL code and query results for calculating the pooled probability.

```
1  WITH total AS(
2    |   SELECT COUNT(DISTINCT uid) AS numcont
3    |   FROM groups
4    ),
5    conversion AS(
6      |   SELECT COUNT(DISTINCT activity.uid)
7      |   FROM activity
8      |   LEFT JOIN groups ON activity.uid=groups.uid
9    ),
10   pooledprob AS (
11     |   SELECT CAST((SELECT * FROM conversion) AS FLOAT)*100 /
12     |   (SELECT * FROM total) AS rate)
13   |   select * from pooledprob
```

Figure 10 SQL code for pool probability.

Query Results	
1 ROWS	
rate	FLOAT8
4.2784463559651025	

Figure 9 Pooled probability result.

Pooled Standard Error

The pooled standard error which is given by this formula:

$$SE_{pool} = \sqrt{P_{pool} * (1 - P_{pool}) * \left(\frac{1}{Total_A} + \frac{1}{Total_B}\right)}$$

SQL code & Results

The following is the SQL code and query results for calculating the pooled standard error.

```
1  with control as(
2    select count(distinct uid) as totalA
3    from groups as g
4    where g.group='A'
5  ),
6  numactA as(
7    select count(distinct activity.uid)
8    from activity
9    left join groups
10   on activity.uid=groups.uid
11   where groups.group = 'A'),
12   treat as(
13    select count(distinct uid) as totalB
14    from groups as g
15    where g.group='B'),
16  numactB as(
17    select count(distinct activity.uid)
18    from activity
19    left join groups
20   on activity.uid=groups.uid
21   where groups.group = 'B'),
22  ppool AS (
23    SELECT CAST((SELECT * FROM numactA)+(SELECT * FROM numactB) AS FLOAT) /
24    | ((SELECT * FROM control)+(SELECT * FROM treat)) AS rate,
25    SELECT SQRT(rate*(1-rate)*((1.0/totalA)+(1.0/totalb))) AS se
26    FROM ppool, control,treat
```

Figure 11 Pooled standard error code.

Query Results	
1 ROWS	
se	FLOAT8
0.001829526081285274	

12 Calculated pooled standard error.

Hypothesis Test

To assess the impact of the new version, a hypothesis test was conducted to determine if there is a significant difference between the conversion rates of the two groups. This test plays a crucial role in informing our decision on whether to proceed with the launch or not. If a noticeable difference is observed in the conversion rates, it provides valuable evidence to support the decision-making process. On the other hand, if the difference is found to be statistically insignificant, it suggests that further investigation and additional data are necessary to arrive at a conclusive decision. In such cases, it is advisable to revisit the problem and gather more data to make a well-informed conclusion. The hypothesis test serves as a critical tool in evaluating the effectiveness of the new version and guiding the decision-making process for the future direction of the project.

$$\hat{d} = \text{Conversion Rate}_B - \text{Conversion Rate}_A$$

Hypothesis:

$$H_0: \hat{d} = 0, \quad \hat{d} \sim N(0, SE_{pool})$$

$$H_a: \hat{d} \neq 0$$

If the following condition is true, then reject the null hypothesis. If not, then we cannot conclude anything.

$$\hat{d} > 1.96 * SE_{pool} \text{ or } \hat{d} < -1.96 * SE_{pool}$$

$$1.96 * SE_{pool} = 0.003528$$

Difference Calculation

In the previous calculations we had the following results:

$$\text{Conversion rate}_B = 4.630$$

$$\text{Conversion rate}_A = 3.923$$

$$\hat{d} = 4.630 - 3.923 = 0.707$$

Data Visualization

I used Jupyter notebooks and Pandas to connect to the database and acquire more information about the data. The Jupyter notebook code is submitted in the zipped file.

Firstly, I plotted the normal distribution curves for the 3 observations (group A and group B and difference between groups' conversion rates) using their standard deviations and means. As observed, changing between the two groups resulted in a normal distribution curve that does not have a mean of 0.

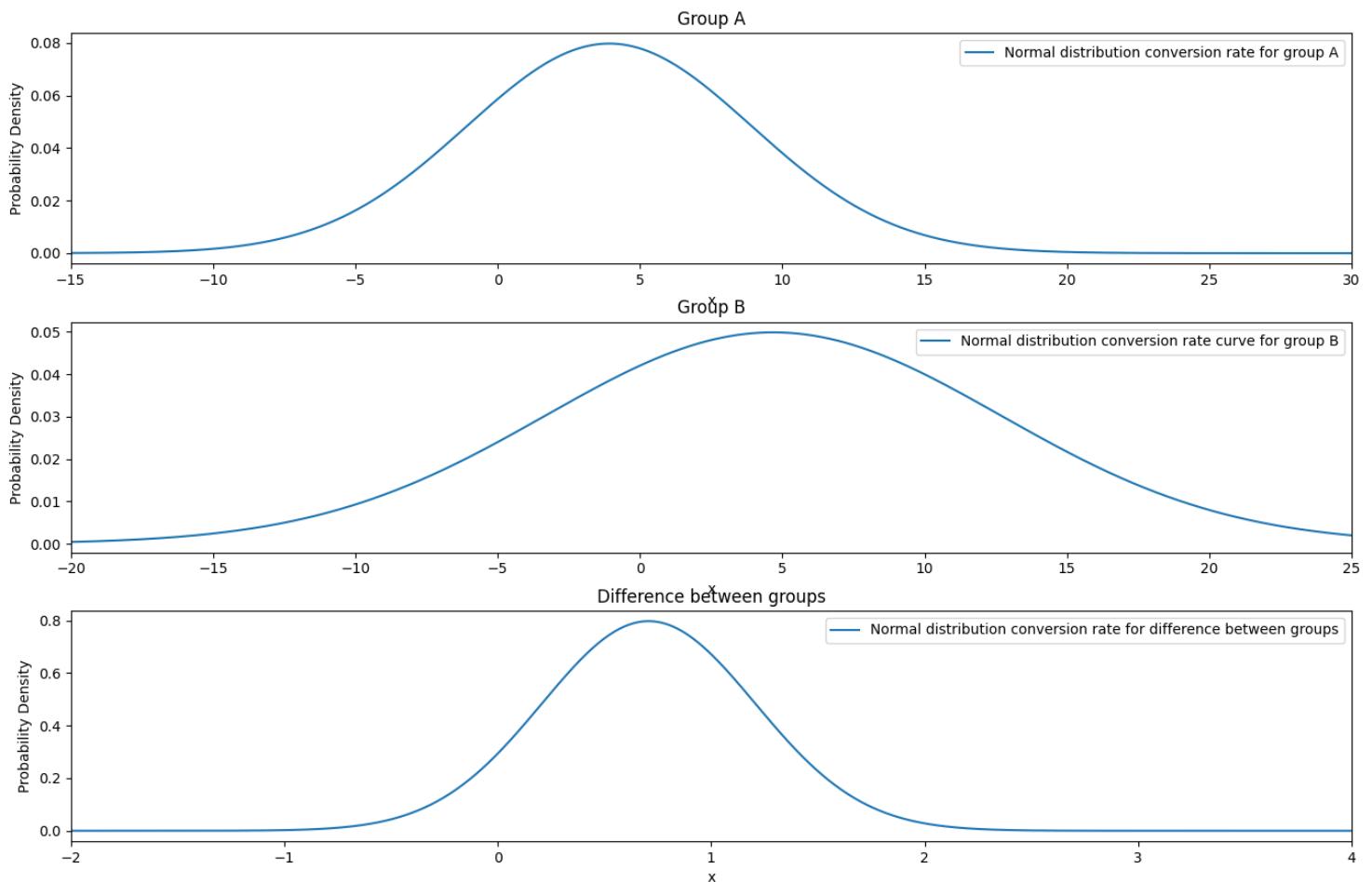
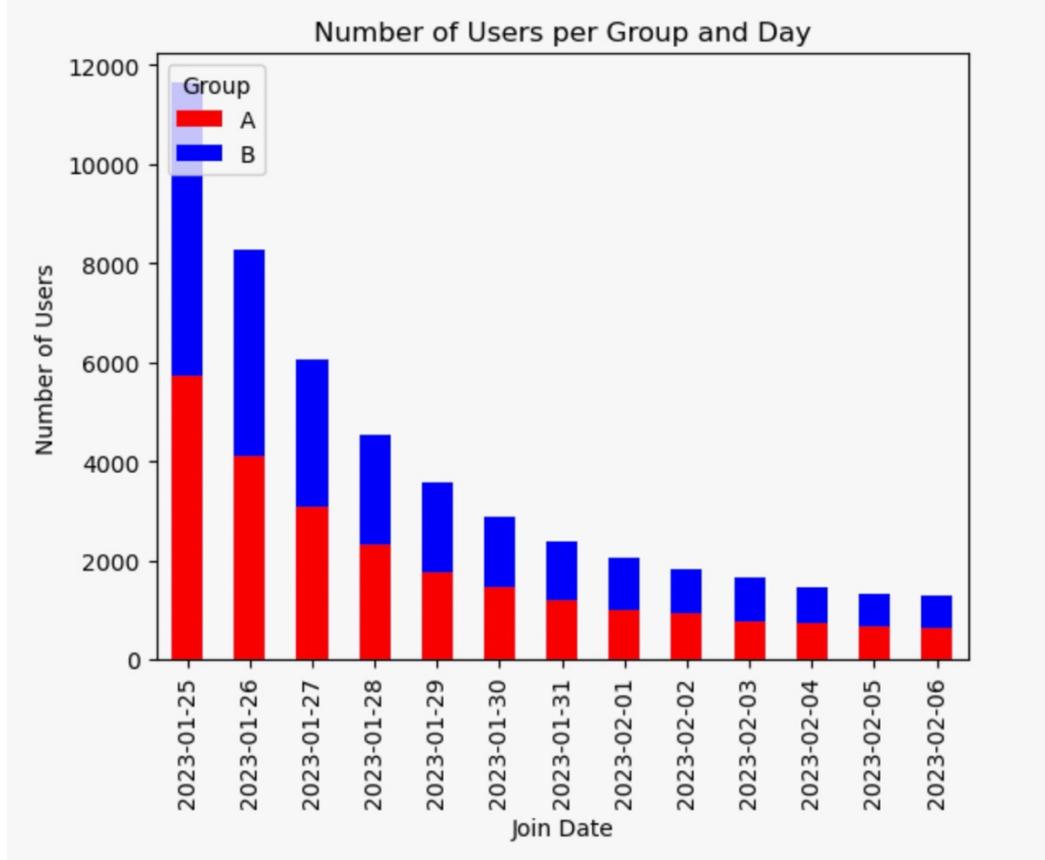
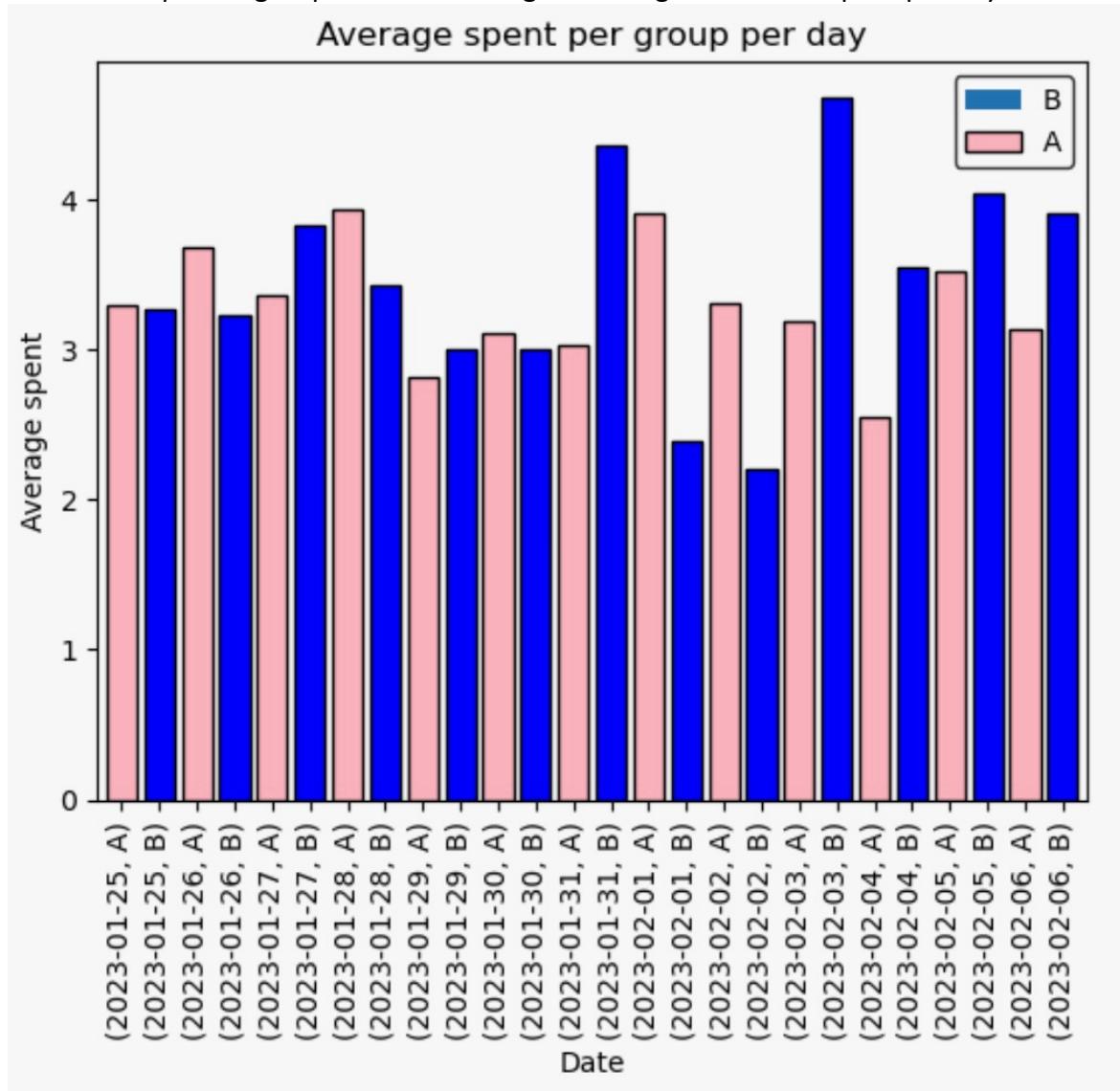


Figure 13 Normal Distribution Curves.

Then, I plotted the number of visitors to the website on each day for each group. Which is shown in Figure 14 below. By analysing this Figure, we can see that the number of visitors per day have increased for Group B. Which means that group B has more visitors and therefore it would be better to launch the new version.



Finally, I used the simplest factor of all “average amount spend per day for each group”, and as expected group B has overall higher average of amount spent per day.



Conclusion & Recommendations

The calculated difference (\hat{d}) between the conversion rates of the two groups, which is found to be greater than the value of 0.003528, indicates a statistically significant distinction. Consequently, we reject the null hypothesis that assumes the equality of the conversion rates. Moreover, this result provides compelling evidence to support the decision to proceed with the launch of the new version. By demonstrating a notable difference between the groups, the analysis reinforces the rationale behind moving forward with the implementation of the updated version, thereby reinforcing the potential benefits it may bring. Furthermore, the Pandas visualizations for different factor also were in favor of launching the new version. So based on our statistical analysis we can confidently conclude that the new version should be launched.