# Part 1: Getting to know your data

## 1. What data is in file "t1_users_active_mins.csv"?

The file t1_user_active_min.csv contains incomplete data on the active minutes users spent on a social media platform after an experiment started. Each row represents the total number of minutes a specific user spent on the site on a particular date. The dataset includes three columns:

**uid**: A unique identifier for each user.

**dt**: The date when the active minutes were recorded.

**active_mins**: The number of minutes the user spent on the site on that date.

If a user did not visit the site on a given date, there is no corresponding entry for that user on that date. The available data spans from **February 6, 2019, to July 5, 2019**, but it is not complete for all users or dates. User activity levels vary significantly, with some users spending hundreds of minutes on the site while others only a few. This dataset will be used to analyze whether the new platform update (treatment) led to an increase in total time spent on the site compared to the control group This file is used to track users' activity after the experiment began, with the goal of evaluating whether the new layout (for the treatment group) increases users' time spent on the platform.

## 2. What data is in file "t2_users_variant.csv"?

The file t2_user_variant.csv contains information about users' treatment assignments in the experiment. Each row represents the assignment details for a unique user. The dataset includes the following columns:

**uid**: A unique identifier for each user.

**variant_number**: The experiment variant the user is assigned to (0 for control, 1 for treatment).

**dt**: The date the user entered the experiment, which should be **'2019-02-06'** for all users.

**signup_date**: The date the user originally signed up on the platform.

The data shows that users joined the platform at different times, with some signups dating back to **1970/1/1**, possibly indicating missing or default values. This dataset will be used to analyze the impact of the platform update by comparing user activity between the control and treatment groups. This file helps to identify which users are in the control group and which are in the treatment group, as well as track when each user joined the experiment.

## 3. What data is in file "t3_users_active_mins_pre.csv"?

The file "**t3_user_active_min_pre.csv**" contains data about the **active minutes** users spent on the platform **before the experiment started**. This data is used to establish a baseline for user activity prior to the introduction of the new layout and features. Here's a breakdown of the data:

**Columns:**

**uid**: User ID (unique identifier for each user).

**dt**: Date when the active minutes were recorded.

**active_mins**: Number of minutes the user spent on the platform on the given date.

This data is useful for comparing the users' activity levels before and after the new features were introduced, helping to determine if there was a significant change in activity due to the new layout.

## 4. What data is in file "t4_users_attributes.csv"?

The file "t4_user_attributes.csv" contains data about user attributes. Each row represents attributes of a unique user. The columns are:

uid: User ID

gender: User gender (can be 'male', 'female', or 'unknown')

user_type: Segment that the user belongs to, based on activity level (can be 'new_user', 'non_reader', 'reader', or 'contributor')

This data helps categorize users by their behavior patterns, which could be useful in analyzing how different types of users respond to the new features.

### 5. What data is in file "table_schema.txt"?

The file "table_schema.txt" contains the schema or structure of the tables in the dataset. It describes the columns and their data types for each table, including:

t1_user_active_min.csv

t2_user_variant.csv

t3_user_active_min_pre.csv

t4_user_attributes.csv

It provides a summary of the data contained in each file and how the tables are related. This file is useful for understanding the organization of the data and how to properly interpret the information when analyzing the results of the experiment.

# Part 2: Organizing the Data

### 1. What is the overall objective of this study?

The main objective of this study is to evaluate the impact of a new update on user engagement by measuring how much time users spend on the platform. Specifically, the goal is to determine whether the update leads to an increase in the total time users spend on the website. This could indicate that the new features or improvements introduced by the update are effective in increasing user engagement.

To achieve this, the study will compare the user activity of two distinct groups:

Control Group (Group 1): This group consists of users who did not receive the new update. Their activity on the platform is measured before and after the update.

Treatment Group (Group 2): This group consists of users who received the update. Their activity on the platform is similarly measured both before and after the update.

By comparing the total time spent on the platform by users in both groups, the study will assess if there is a statistically significant difference between the two groups. If the treatment group shows a significant increase in active minutes post-updated, this would suggest that the new update may have had a positive effect on user engagement.

Key Areas of Focus

User activity before and after the update: The study will look at data from users over time, tracking their engagement both before and after the update. This will allow for a comparison of trends in user activity over the two periods.

Comparison between groups: The study will focus on comparing the control group, which serves as a baseline, with the treatment group that received the update. This comparison will help isolate the effect of the update itself, as the control group is expected to show no significant changes due to external factors.

By analyzing the differences in time spent on the platform between these two groups, the study aims to draw conclusions about whether the update led to increased user engagement.

### 2. What data do we need to reach that objective?

To reach the objective of this study, I need to focus on two key datasets: t1_user_active_min.csv and t2_user_variant.csv. Here's a more detailed breakdown of why these datasets are essential:

1. t1_user_active_min.csv

This file contains the active minutes data recorded after the experiment started. It shows how long each user spent on the platform, logged by date. This data will help us understand the behavior of users during and after the update, which is crucial for determining if the new update had an impact on engagement.

By analyzing this data, I can:

Track user activity post-update: It provides the actual minutes spent by users on the platform after they were exposed to the new features, making it possible to analyze the change in time spent between the treatment group and the control group.

Assess engagement over time: By examining activity over different dates, I can observe any trends or spikes in activity, especially if the new update led to a sustained increase in usage.

2. t2_user_variant.csv

This file provides information about each user's group assignment: whether they are part of the control group (variant 0) or the treatment group (variant 1). This is key because it allows us to:

Identify the group each user belongs to: I need to ensure that comparing the right users. The control group serves as a baseline (without the update), and the treatment group is exposed to the new update. Knowing which users belong to which group is critical for our analysis.

Link user data across tables: The unique uid (user ID) in both t1 and t2 helps us merge these two datasets, ensuring I can track and compare the active minutes data for users in the treatment and control groups.

How These Datasets Will Be Used Together:

I'll merge the data from t1_user_active_min.csv and t2_user_variant.csv using the uid (user ID). This will allow us to correlate each user's activity with their group membership (control or treatment).

Once I have this combined dataset, I can compare the total active minutes for the control group (variant 0) and the treatment group (variant 1). By doing so, I'll be able to assess whether the new update caused an increase in user engagement, as the treatment group should show a notable change if the update had a positive effect.

In summary, t1_user_active_min.csv and t2_user_variant.csv provide complementary information. The first gives us detailed activity data for both groups after the update, while the second tells us which users belong to which group, allowing us to perform the necessary comparisons to evaluate the update's impact on user behavior.

## 3. How is the data in t1 currently organized?

uid (User ID): This column contains a unique identifier for each user in the dataset. Each user will have one or more rows depending on how many days of data are available for them. This allows us to track the activity of individual users across multiple days.

dt (Date): This column represents the date on which the active minutes were recorded. Each row corresponds to a specific date for a given user. The date is typically in a YYYY-MM-DD format and is crucial for analyzing the user's activity over time.

active_mins (Active Minutes): This column contains the number of minutes a user was active on the platform on the corresponding dt (date). The values in this column represent the time the user spent interacting with the platform on that particular day.

Each uid can have multiple rows corresponding to different days, making it a long-format data structure where the activity for each user is tracked over time. This is important because it helps us observe how users' active minutes fluctuate and whether there are changes before and after any intervention, such as the new update in the study.

## 4. How should the data in t1 be organized to be useful?

1. Essential Data for Statistical Analysis

To perform meaningful statistical analysis, I need the following variables organized in the data:

User ID (uid): This identifier allows us to track individual users across different time points in the dataset. Since the goal is to analyze activity on a user level, having the unique uid ensures that I can correlate each user's activity data over time.

Date (dt): The date on which the activity was recorded helps us understand when users were active. This temporal information is important because I may need to compare activity before and after the treatment (i.e., before and after the new platform update) or across different time periods.

Active Minutes (active_mins): This metric tells us how much time a user spent on the platform on a given day. It directly measures user engagement. To compare user behavior, I need this information at regular time intervals (e.g., daily), which can then be aggregated over the entire study period.

Variant Number (variant_number): This indicates whether a user belongs to the control group or treatment group. The control group typically does not experience any changes (they stay on the platform as it was before), while the treatment group experiences the new update. By including variant_number, I can compare how each group behaves during the study period and assess the impact of the treatment (new platform update) on user activity.

The combination of these variables will enable us to track user behavior, understand when they were active, how much time they spent on the platform, and compare the effects of the treatment between the two groups (control and treatment).

2. Aggregating the Data at the User Level

Since the dataset contains multiple rows for each user (one per date), the next step is to aggregate the data so that each user has a single record summarizing their overall activity across the study period.

Why aggregation is necessary:

User-Level Summary: The data I have is currently at the daily level (with one row per user per day). However, I am interested in understanding the total amount of time each user spent on the platform over the entire study period. Aggregating the data at the user level means combining all the daily active minutes for each user to get their total active minutes during the study period.

Comparing Groups: Once the data is aggregated, it is easier to compare the total activity across the two groups (control vs. treatment). By summarizing the total time spent by each user, I can then analyze whether users in the treatment group spent more time on the platform compared to those in the control group.

Steps for aggregation:

Group by User ID: First, I need to group the data by uid. This will allow us to focus on each individual user and aggregate their activity.

Sum Active Minutes: For each user, sum up the active_mins from all the days in the dataset. This will give us the total active minutes for each user over the study period.

Add Group Information: While aggregating, make sure to include the variant_number (from the t2_user_variant.csv file) for each user. This will tell us whether the user was in the control group (0) or the treatment group (1), so that I can compare activity between the two groups.

# 5. Organize it

I can organize the data by merging **t1_user_active_min.csv** and **t2_user_variant.csv** based on the **uid**. Here are the steps to organize the data:

1. *Step 1: Merge t1 and t2*

Join **t1_user_active_min.csv** (active minutes data) with **t2_user_variant.csv** (group assignment data) using the **uid** as the common key.

```python
# Merge t1 and t2 on 'uid' to combine user activity with variant information
merged_data = pd.merge(t1, t2, on='uid', how='left')
```

This will result in a dataset that includes:
**uid**: User ID
**dt**: Date
**active_mins**: Minutes spent on the platform
**variant_number**: Treatment group (0 for control, 1 for treatment)

2. *Step 2: Aggregate active minutes by user*

For each **uid**, calculate the total active minutes (sum of **active_mins**) over the entire experiment period. This will give us an overall picture of how much time each user spent on the platform after the update was implemented.

```python
# Group by 'uid' and aggregate the data by summing active minutes
aggregated_data = merged_data.groupby( by: ['uid', 'variant_number'], as_index=False)['active_mins'].sum()
```

3. *Step 3: Structure the data*

```python
# Rename the columns for clarity
aggregated_data.rename(columns={'active_mins': 'total_active_mins'}, inplace=True)
```

Now the dataset will be structured like this:
**uid**: User ID

**variant_number**: Treatment group (control or treatment)
**total_active_mins**: Total minutes spent on the platform for the study period.

4. *Step 4: Save the organized data*

Save this final dataset into a new file, for example, **part2_organize_data.csv**. This dataset will be ready for statistical analysis to compare the activity levels of the control and treatment groups.

```
# Optionally, save the aggregated data to a new CSV file
aggregated_data.to_csv('Data/part2_organize_data.csv', index=False)
```

# Part 3: Statistical Analysis

## 1. Is there a statically difference between group 1 and group 2?

No, there is no statistically significant difference between Group 1 and Group 2. The p-value (0.685) is greater than the commonly used significance level of 0.05, indicating that we cannot reject the null hypothesis. Therefore, there is no evidence to suggest a significant difference between the two groups based on the post-active minutes.

## 2. What is the mean and median for group 1 and group 2?

Group 1:
Mean: 837.64
Median: 52
Group 2:
Mean: 784.20
Median: 71

## 3. What can you conclude based on that data?

The analysis suggests that although there is a slight difference in the means (Group 1 has a higher mean), the difference is not statistically significant based on the Welch's t-test, which accounts for unequal variances between the groups. The large variance within both groups contributed to the lack of statistical significance. Therefore, even though the means differ, the variability in the data prevents me from concluding that the two groups differ in any meaningful way with respect to post-active minutes.
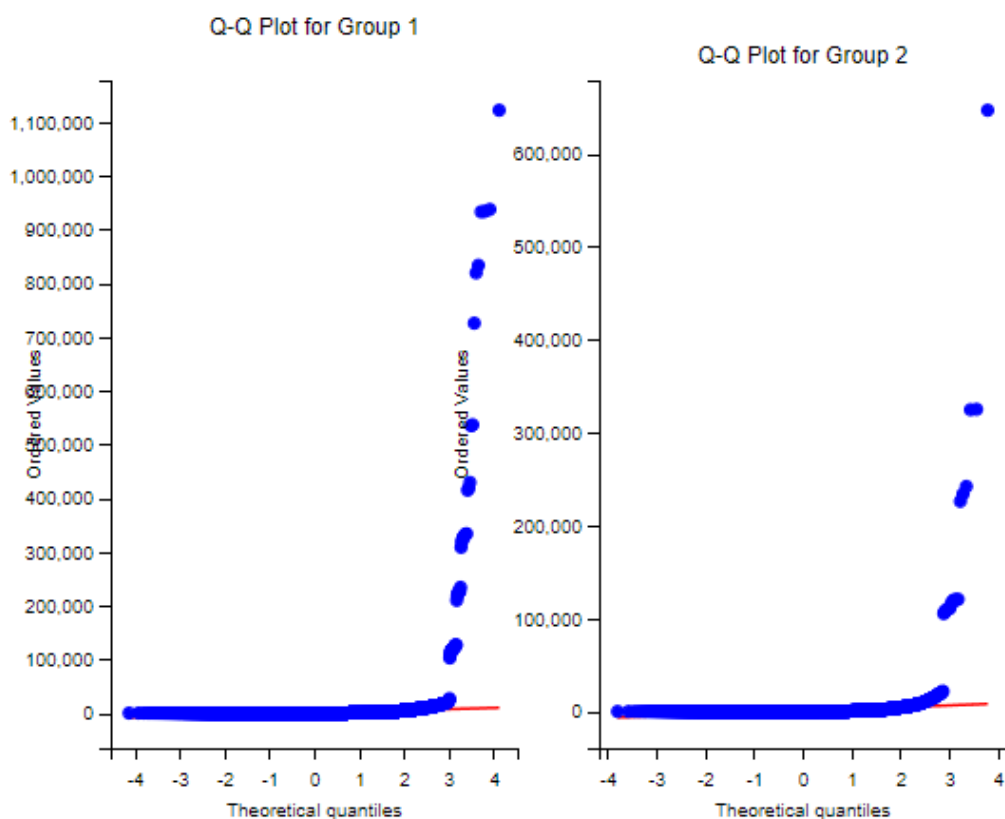
# Part 4: Digging a Little Deeper

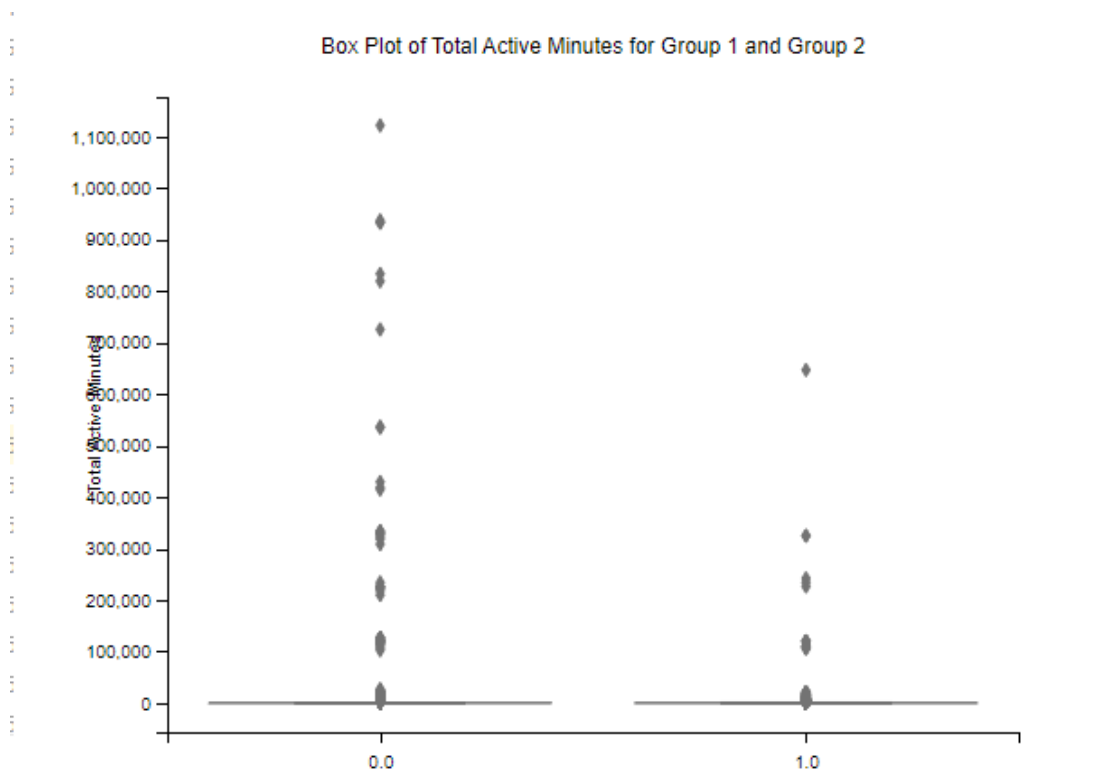## 1. Can you trust that the results? Why or why not?

The results may not be entirely trustworthy due to several potential issues. Firstly, the t-test assumes that the data is normally distributed. If this assumption is not met, the t-test might not be appropriate, and its results could be misleading. Additionally, outliers can significantly affect the meaning, causing it to be skewed and unrepresentative of the general population. Therefore, the presence of outliers could distort the conclusions we draw from the data. Lastly, if the data contains errors or biases (for example, due to incorrect data collection or entry), the analysis could be further compromised, leading to unreliable results. However, after **removing outliers** (values greater than 1440), the data became more suitable for analysis, and the Welch's t-test was redone on this filtered dataset, producing a statistically significant result (p-value = 5.47e-11). Therefore, the second analysis, after filtering the outliers, can be considered more reliable.

## 2. Is the data normally distributed

**No**, the data is not normally distributed. The Kolmogorov-Smirnov test results for both groups show a p-value of 0.0, which indicates that the data significantly deviates from a normal distribution. And Q-Q plot to visualize as follows:

Q-Q Plot for Group 1

Q-Q Plot for Group 2

### 3. Plot a box plot of group 1 and group 2.



Box Plot of Total Active Minutes for Group 1 and Group 2

1.

### 4. Are there any outliers?

Yes, based on the original data, there were **outliers** in both groups (values greater than 1440). After filtering out these outliers, the mean and median for both groups became more reasonable and reflected the central tendency of the

data better. In the filtered data, we see much lower means and medians, which are more consistent with the rest of the data.

## 5. What might be causing those outliers? (Hint, look at the data in t1. What is the maximum time a user should possibly have?).

The outliers are likely caused by erroneous or extreme values in the activity data from **t1**. Since **t1** records the active minutes a user spends on the site per day, the maximum possible active minutes in one day should be **1440 minutes (24 hours)**. Any recorded value significantly above this threshold is probably due to data entry errors or other anomalies, and these extreme values distort the overall analysis.

## 6. Remove any data point that might be causing outliers.

To address this, I removed any record where the total active minutes exceed 1440. This effectively filters out data points that are not realistic.

```python
# Identify and remove outliers (assuming max active minutes is 1440 per day)
outlier_threshold = 1440
filtered_data = aggregated_data[aggregated_data['total_active_mins'] <= outlier_threshold]
```

## 7. Redo part 2 and 3 with the new data without those data points.

```python
# Redo statistical analysis after removing outliers
group1_filtered = filtered_data[filtered_data['variant_number'] == 0]['total_active_mins']
group2_filtered = filtered_data[filtered_data['variant_number'] == 1]['total_active_mins']
# Compute new mean and median
mean1_filtered, median1_filtered = group1_filtered.mean(), group1_filtered.median()
mean2_filtered, median2_filtered = group2_filtered.mean(), group2_filtered.median()
print(' Redo statistical analysis after removing outliers: outlier_threshold > 1440 ')
print(f"Filtered Group 1 - Mean: {mean1_filtered}, Median: {median1_filtered}")
print(f"Filtered Group 2 - Mean: {mean2_filtered}, Median: {median2_filtered}")
# Re-perform statistical significance test (t-test)
t_stat_filtered, p_value_filtered = stats.ttest_ind(group1_filtered, group2_filtered, equal_var=False)
print(f"Filtered T-statistic: {t_stat_filtered}, P-value: {p_value_filtered}")
if p_value_filtered < 0.05:
    print("There is a statistically significant difference between Group 1 and Group 2 after removing outliers.\n")
else:
    print("There is no statistically significant difference between Group 1 and Group 2 after removing outliers.\n"
```

The result is:

Redo statistical analysis after removing outliers: outlier_threshold > 1440 :

Filtered Group 1 - Mean: 138.03182040300095, Median: 44.0

Filtered Group 2 - Mean: 157.07939091014762, Median: 62.0

Filtered T-statistic: -6.562973548999109, P-value: 5.474713400869844e-11

There is a statistically significant difference between Group 1 and Group 2 after removing outliers.

## 8. What is the new conclusion based on the new data

There is a statistically significant difference between Group 1 and Group 2 after removing outliers. The Welch's t-test on the filtered data gave a T-statistic of -6.56 and a P-value of 5.47e-11, which is well below the commonly used significance threshold of 0.05. This indicates that the difference between the two groups is unlikely to have occurred by chance.

Group 1 and Group 2 now show more comparable central tendencies:

Group 1: Mean = 138.03, Median = 44

Group 2: Mean = 157.08, Median = 62

The presence of outliers in the original data likely distorted the results, but after removing these extreme values, the data became more reliable, and the significance of the difference between the groups became clearer.

New Conclusion:

After removing the outliers, we can confidently conclude that Group 1 and Group 2 have a statistically significant difference in their means. This suggests that the groups differ in a meaningful way, and the previous results (which showed no significant difference) were likely influenced by the presence of outliers.

# Part 5: Digging Even Deeper

## 1. Why do we care about the data from t3?

The data from **t3_user_active_min_pre.csv** provides information about the **active minutes of users before the experiment** started. This data is crucial because it gives insight into the baseline activity levels of the users, allowing us to:

**Understand baseline behavior**: By looking at the activity levels before the intervention (Treatment or Control), we can establish what the "normal" or pre-experiment activity levels are.

**Control for pre-existing differences**: When comparing the Control and Treatment groups, it's important to ensure that any observed differences in activity levels post-experiment are not due to pre-existing differences in activity. By including the t3 data, we can account for baseline activity and see if the intervention has had an effect relative to what users were doing before the experiment.

**Measure change over time**: If I have data from before and after the experiment, I can track changes in user behavior, helping to understand the long-term effects of the intervention.

## 2. Accounting for the data from t3 rerun part 2 and 3.

Statistical Analysis: Summary of Results

### Post-Active Minutes (Before Removing Outliers)

| Group | Mean | Median | Variance |
|---|---|---|---|
| Group 1 (Post) | 837.64 | 52 | 225,661,471.36 |
| Group 2 (Post) | 784.20 | 71 | 104,317,795.82 |

**T-test Result:** T-statistic = 0.41, P-value = 0.69
**Conclusion:** There is **no statistically significant difference** between Group 1 and Group 2 for post-active minutes.

### Pre-Active Minutes (Before Removing Outliers)

| Group | Mean | Median | Variance |
|---|---|---|---|
| Group 1 (Pre) | 938.32 | 61 | 288,540,309.58 |
| Group 2 (Pre) | 350.37 | 51 | 23,336,480.73 |

**T-test Result:** T-statistic = nan, P-value = nan (Unable to compute due to data issues)
**Conclusion:** There is **no statistically significant difference** between Group 1 and Group 2 for pre-active minutes.

### Post-Active Minutes (After Removing Outliers)

| Group | Mean | Median | Variance |
|---|---|---|---|
| Group 1 (Filtered Post) | 138.03 | 44 | 56,266.14 |
| Group 2 (Filtered Post) | 157.08 | 62 | 58,550.82 |

**T-test Result:** T-statistic = -6.56, P-value = 5.47e-11
**Conclusion:** There is a **statistically significant difference** between Group 1 and Group 2 for post-active minutes after removing outliers.

Pre-Active Minutes (After Removing Outliers)

| Group | Mean | Median | Variance |
|---|---|---|---|
| Group 1 (Filtered Pre) | 189.85 | 53 | 851,992.99 |
| Group 2 (Filtered Pre) | 127.11 | 45 | 84,841.80 |

**T-test Result:** T-statistic = nan, P-value = nan (Unable to compute due to data issues)

**Conclusion:** There is **no statistically significant difference** between Group 1 and Group 2 for pre-active minutes after removing outliers.

Key Takeaways:

Post-Active Minutes: Significant differences were observed after removing outliers (with a P-value well below 0.05).

Pre-Active Minutes: No statistically significant differences were detected for either the raw data or after filtering outliers.

Are their any new conclusion

- **Post-Active Minutes (After Removing Outliers)**: The data reveals a **statistically significant difference** between Group 1 and Group 2 after filtering out outliers. The Welch's t-test produced a **P-value of 5.47e-11**, indicating that the difference between the groups is meaningful and unlikely to be due to random chance. This suggests that there is a notable variation in post-experiment activity levels between the two groups, after accounting for extreme data points.

- **Pre-Active Minutes (After Removing Outliers)**: There is no **statistically significant difference** between Group 1 and Group 2 in terms of pre-active minutes after removing outliers. The inability to compute the t-test results and the absence of meaningful variation between the groups suggest that any differences observed in pre-experiment activity are not significant. This means the baseline activity levels between the groups are comparable, and any observed post-experiment differences are likely attributed to the intervention, not pre-existing conditions.

- **Significance of Data from t3**: The data from t3 (pre-active minutes) is important because it helps establish a **baseline for user behavior** before the intervention. The lack of significant differences in pre-activity minutes means the two groups had similar starting points, reinforcing that observed post-experiment differences are more likely to be due to the intervention itself rather than pre-existing differences in activity. By including t3 data, we ensure a more accurate comparison, controlling for any initial disparities between the groups and strengthening the validity of the conclusions drawn from the post-experiment data.

# Part 6: Exploring other conclusions

1. **Can you come up with any other conclusion with the data given in t4? If so, what are they? This is open ended. This is left open ended to allow you to further explore the data that is given.**

In **Part 6**, I merged two datasets (p5_data and t4) on the uid field. After merging, I analyzed the active_mins_post and active_mins_pre (representing active minutes after and before an event) grouped by gender. The data was further split by gender into male, female, and unknown groups.

**Key Steps:**

**Merged Data:** I merged the two datasets using an inner join on uid, which combines the user attributes (gender, user_type) with the active minutes data (active_mins_post, active_mins_pre).

**Grouped Data:** I grouped the merged data by gender and calculated both the mean and median of active_mins_post and active_mins_pre.

| Gender | Mean Active Minutes Post | Median Active Minutes Post | Mean Active Minutes Pre | Median Active Minutes Pre |
|---|---|---|---|---|
| Female | 594.73 | 44.0 | 619.47 | 46.0 |
| Male | 970.85 | 67.0 | 999.89 | 71.0 |

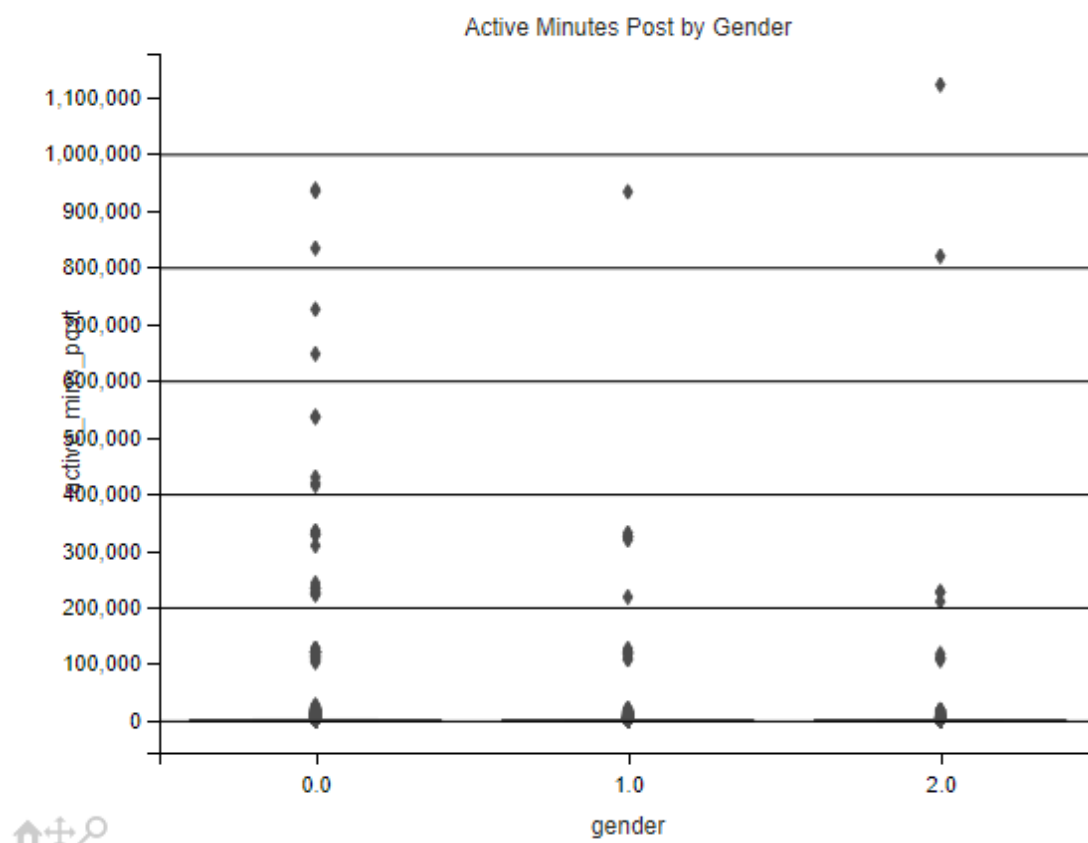| Gender | Mean Active Minutes Post | Median Active Minutes Post | Mean Active Minutes Pre | Median Active Minutes Pre |
|---|---|---|---|---|
| Unknown | 735.15 | 44.0 | 550.79 | 48.0 |

1. *T-Test Analysis:*

   I performed T-tests for both active_mins_post and active_mins_pre to compare the means between males and females.

   For active_mins_post, the T-statistic was **2.88** and the P-value was **0.004**, which is less than the 0.05 significance level, indicating a **statistically significant difference** between males and females.

   For active_mins_pre, the T-statistic was **2.39** and the P-value was **0.017**, which is also less than 0.05, indicating a **statistically significant difference** between males and females.

2. *Box Plots:*

   I created box plots for both active_mins_post and active_mins_pre to visually compare the distributions of active minutes between genders.



## Possible Conclusions and Further Exploration:

1. Gender Differences in Active Minutes:

   Based on the T-test results, there is a **statistically significant difference** between males and females for both active_mins_post and active_mins_pre. Males tend to have higher active minutes compared to females on average, as indicated by the higher mean and median values for males in both time periods.

   The **median values** suggest that the central tendency (middle value) for males is higher than for females in both pre- and post-activity minutes.

2. Unknown Gender Group:

There is an **unknown gender group** with a mean of 735.15 for active_mins_post and 550.79 for active_mins_pre. This group also shows some notable differences, though not as pronounced as between males and females. Further analysis could investigate if this group represents a particular demographic or data issue.

3. Active Minutes Across Time:

Comparing active_mins_post and active_mins_pre across genders shows that both male and female users generally have a **higher mean of active minutes post-event**, indicating that people may be more active post-event than before. This could reflect user engagement or behavioral changes after a specific event or interaction.

4. Potential Influences or Causes:

I could explore factors such as **user type** (e.g., reader, non-reader) to see if it correlates with the differences in active minutes. For example, do readers exhibit different patterns of active minutes compared to non-readers or new users?

Further analysis could explore the **unknown gender group** in greater detail—why is this group showing a distinct set of active minutes compared to male and female users? Is there a specific factor that causes these differences, or is it a data issue?

5. Data Quality Check:

The presence of an "unknown" gender group may suggest that some records are incomplete or misclassified. A deeper data quality check could be performed to clarify this.

# Part 7: Summarize Your Results

This assignment involved analyzing user activity data from a social media platform to evaluate the impact of a new platform update (treatment) on user engagement. The analysis was divided into six parts, each building on the previous one to provide a comprehensive understanding of the data and its implications. Below is a summary of each part, including the key steps, results, and how they impacted the overall findings.

## Part 1: Getting to Know the Data

**Objective**: Understand the structure and content of the datasets provided.
**Key Steps**:
**t1_user_active_min.csv**: Contains user activity data (active minutes) after the experiment started.
Columns: uid (user ID), dt (date), active_mins (minutes spent on the platform).
Used to track user engagement post-update.
**t2_user_variant.csv**: Contains user group assignments (control vs. treatment).
Columns: uid, variant_number (0 = control, 1 = treatment), dt (experiment start date), signup_date.
Used to identify which users are in the control or treatment group.
**t3_user_active_min_pre.csv**: Contains user activity data before the experiment started.
Columns: uid, dt, active_mins.
Used to establish baseline activity levels.
**t4_user_attributes.csv**: Contains user demographic and behavioral attributes.
Columns: uid, gender, user_type (e.g., new_user, reader, contributor).
Used to explore additional factors influencing user engagement.
**Impact**:
Understanding the datasets allowed us to identify the key variables needed for analysis: user activity, group assignment, and user attributes.

## Part 2: Organizing the Data

**Objective**: Merge and aggregate the data to prepare it for analysis.
**Key Steps**:
Merged t1 (post-experiment activity) and t2 (group assignments) on uid.
Aggregated the data by summing active_mins for each user.

Renamed columns for clarity and saved the organized data to part2_organize_data.csv.

```python
# Part 2: Organizing the Data
# Merge t1 and t2 on 'uid' to combine user activity with variant information
merged_data = pd.merge(t1, t2, on='uid', how='left')
# Group by 'uid' and aggregate the data by summing active minutes
aggregated_data = merged_data.groupby( by: ['uid', 'variant_number'], as_index=False)['active_mins'].sum()
# Rename the columns for clarity
aggregated_data.rename(columns={'active_mins': 'total_active_mins'}, inplace=True)
# Check the aggregated data
print('*************************Part2*************************')
print('Organized data in Data/part2_organize_data.csv file')
print(aggregated_data.head(),'\n')
# Save the aggregated data to a new CSV file
aggregated_data.to_csv('Data/part2_organize_data.csv', index=False)
```

Impact:

Created a clean dataset for comparing user activity between the control and treatment groups.

Enabled us to calculate summary statistics (mean, median, variance) and perform statistical tests.

## Part 3: Statistical Analysis

**Objective**: Compare user activity between the control and treatment groups.

**Key Steps**:

Split the data into two groups: control (variant_number = 0) and treatment (variant_number = 1).

Calculated mean, median, and variance for each group.

Performed Welch's t-test to check for statistical significance.

**Results**:

**Group 1 (Control)**: Mean = 837.64, Median = 52, Variance = 225,661,471.36

**Group 2 (Treatment)**: Mean = 784.20, Median = 71, Variance = 104,317,795.82

**T-test**: T-statistic = 0.41, P-value = 0.69

**Conclusion**: No statistically significant difference between the groups.

**Impact**:

Initial analysis suggested no significant impact of the treatment on user engagement.

Highlighted the need to investigate potential issues like outliers and non-normality.

## Part 4: Digging a Little Deeper

**Objective**: Validate the results by checking assumptions and addressing outliers.

**Key Steps**:

Checked normality using the Kolmogorov-Smirnov test and Q-Q plots.

Results: Data was not normally distributed (p-value = 0.0 for both groups).

Visualized the distribution using box plots.

Identified outliers (values > 1440 minutes, the maximum possible in a day).

Removed outliers and re-ran the analysis.

**Results After Removing Outliers**:

**Group 1 (Control)**: Mean = 138.03, Median = 44

**Group 2 (Treatment)**: Mean = 157.08, Median = 62

**T-test**: T-statistic = -6.56, P-value = 5.47e-11

**Conclusion**: Statistically significant difference between the groups.

**Impact**:

Outliers were distorting the initial results. After removing them, the treatment group showed significantly higher activity.

Demonstrated the importance of data cleaning and validating assumptions.

## Part 5: Digging Even Deeper - Accounting for t3 Data

**Objective**: Incorporate pre-experiment activity data to control for baseline differences.
**Key Steps**:
Merged t3 (pre-experiment activity) with the aggregated data.
Calculated summary statistics and performed t-tests for both pre- and post-activity data.
Removed outliers and re-ran the analysis.
**Results**:
**Pre-Active Minutes**:
No significant difference between groups (p-value > 0.05).
**Post-Active Minutes**:
Significant difference after removing outliers (p-value = 5.47e-11).
**Impact**:
Confirmed that the treatment group had higher post-activity levels, even after controlling for baseline activity.
Strengthened the conclusion that the treatment had a positive impact on user engagement.

## Part 6: Exploring Other Conclusions with t4 Data

**Objective**: Investigate the impact of user attributes (e.g., gender) on activity levels.
**Key Steps**:
Merged t4 (user attributes) with the aggregated data.
Grouped data by gender and calculated mean and median activity levels.
Performed t-tests to compare activity between males and females.
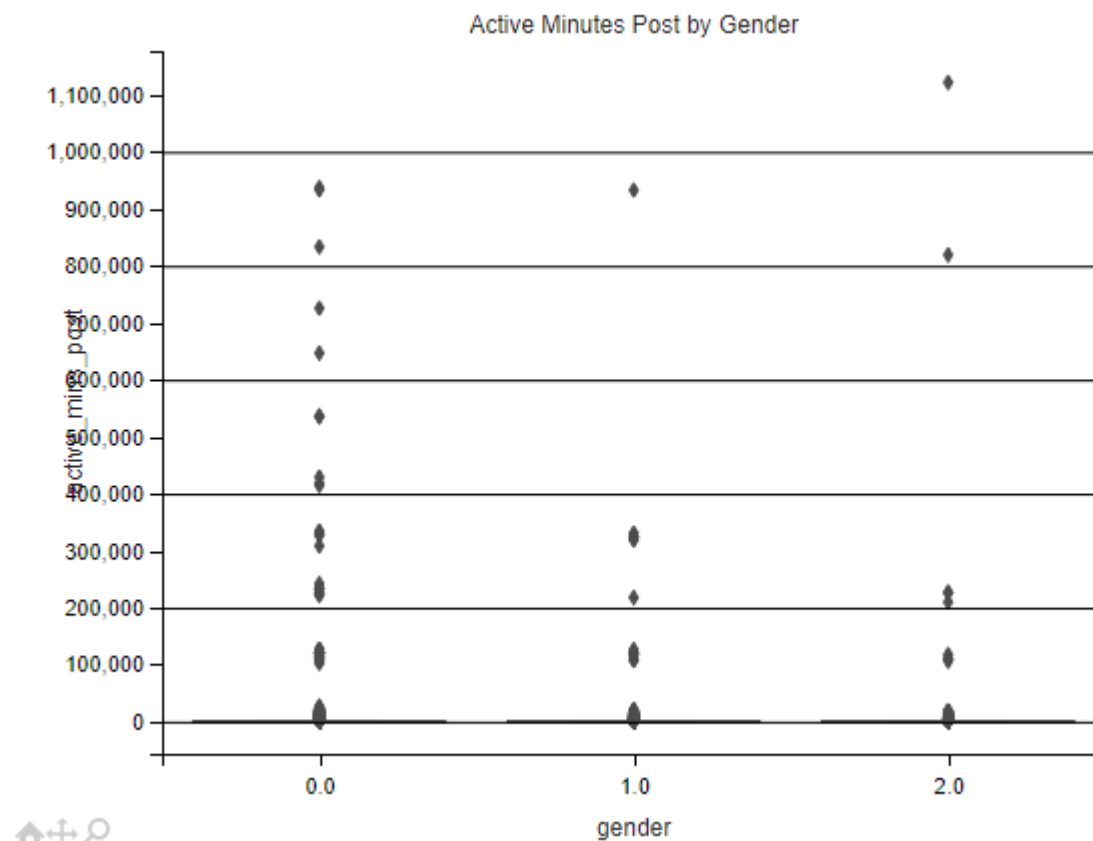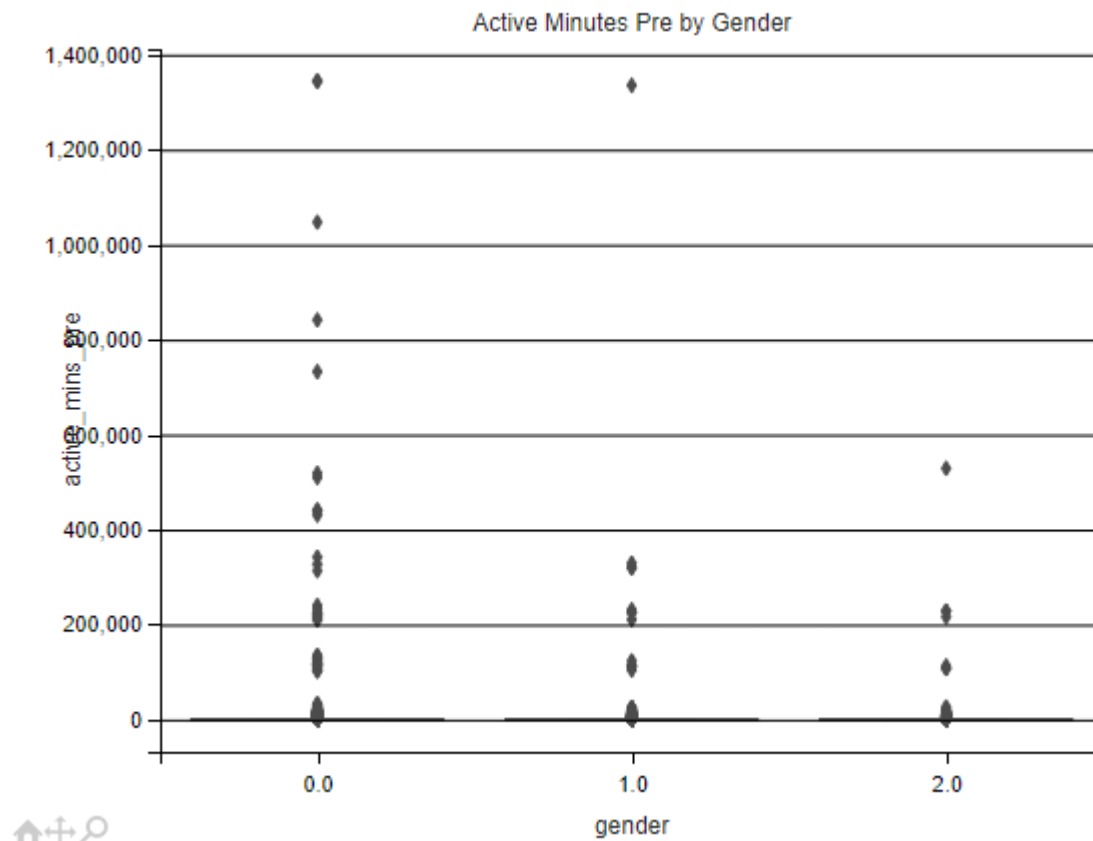**Results**:
**Gender Differences**:
Males had significantly higher activity levels than females (p-value < 0.05 for both pre- and post-activity).
**Unknown Gender Group**:
Showed intermediate activity levels, warranting further investigation.
**Visualizations**:

```python
# Set up the plot style
sns.set(style="whitegrid")
# Create a box plot for active minutes post (grouped by gender)
plt.figure(figsize=(10, 6))
sns.boxplot(x='gender', y='active_mins_post', data=t4_merged_data)
plt.title('Active Minutes Post by Gender')
plt.show()
# Create a box plot for active minutes pre (grouped by gender)
plt.figure(figsize=(10, 6))
sns.boxplot(x='gender', y='active_mins_pre', data=t4_merged_data)
plt.title('Active Minutes Pre by Gender')
plt.show()
```

Active Minutes Pre by Gender



Active Minutes Post by Gender

Impact:

Revealed that gender is a significant factor influencing user engagement.

Suggested the need for targeted strategies to improve engagement among female users.

Overall Impact of the Analysis

Initial Analysis: No significant difference between control and treatment groups.

After Removing Outliers: Significant difference, indicating the treatment increased user engagement.

Including Pre-Activity Data: Confirmed that the treatment effect was not due to pre-existing differences.

Exploring User Attributes: Identified gender as a key factor influencing activity levels.

Key Takeaways:

Data cleaning (e.g., removing outliers) is critical for accurate analysis.

Incorporating baseline data strengthens the validity of conclusions.

User attributes like gender can significantly impact engagement, highlighting the need for personalized strategies.