

A Brief Description Of The Dataset And A Summary Of Its Attributes

The dataset contains information about individual rides in a bike-sharing system covering the greater San Francisco Bay area. It includes details such as start and end stations, start and end times, user type (customer or subscriber), member gender, and year of birth. The dataset has undergone preliminary wrangling, including the removal of missing data and conversion of the member birth year to integers. Exploratory data analysis (EDA) has been conducted to uncover trends and insights.

Initial Plan For Data Exploration

The initial plan for data exploration involved investigating trends in user age, user status (subscriber or customer), gender distribution, station popularity, and user activity times. The main features of interest were identified, and the dataset was cleaned and transformed accordingly to support the analysis.

Actions Taken For Data Cleaning And Feature Engineering

- Data cleaning involved dropping rows with missing values.
- Converting the member birth year to integers, and extracting additional features such as year, month, day, and user age.
- Feature engineering included converting the duration from seconds to minutes and selecting relevant columns for further analysis.

Key Findings And Insights

- Subscribers make up the majority of users (90.5%), and there are more male users (74.6%) than any other gender.
- Most users prefer not to share bikes for all trips (over 90%).
- Thursdays are the busiest days for both trips starts and ends.
- Subscribers tend to be older than customers, and the majority of male users are subscribers.
- The duration of rides is influenced by user type, with subscribers having shorter rides on average.

Formulating At Least 3 Hypotheses About This Data

Age and Subscription Type: Subscribers are likely to be older than customers.

Gender and Bike Sharing: Female users are more likely to share bikes for all trips compared to male users.

Day of the Week and Ride Duration: The duration of rides is longer on weekends compared to weekdays.

Conducting A Formal Significance Test For One Of The Hypotheses And Discussing The Results

For hypothesis 1, a t-test could be performed to compare the mean ages of subscribers and customers. The null hypothesis would be that there is no significant difference in the mean ages of the two groups. The alternative hypothesis would be that subscribers are older than customers.

Conducting the t-test and analyzing the p-value would help determine whether to reject the null hypothesis.

Suggestions For The Next Steps In Analysing This Data

Further Investigation of Age Groups: Explore age groups to identify specific age ranges with different usage patterns.

Temporal Patterns: Investigate temporal patterns, such as daily and monthly variations in bike usage.

Geospatial Analysis: Explore the geographical distribution of bike stations and user activity.

Machine Learning Models: Consider building predictive models to forecast bike usage based on various features.

Summary Of The Quality Of This Dataset And A Request For Additional Data If Needed

The dataset appears to be of good quality after the cleaning and wrangling process. However, additional data on user demographics, station locations, and more detailed trip information could enhance the depth of analysis and provide more comprehensive insights into user behaviour and system usage.