

# ***A New Hint to Transportation Analysis of the NYC Bike Share System***

***A Project Reporting***

***Submitted by***

**A.AJAY**

**510519104002**

**R.MADHAN**

**510519104015**

**K.THIRUPATHI**

**510519104029**

**P.VIGNESH**

**510519104033**

***in partial fulfillment for the reward of the degree***

***of***

**BACHELOR OF ENGINEERING**

**IN**

**COMPUTER SCIENCE AND ENGINEERING**

**BHARATHIDASAN ENGINEERING COLLEGE**



**NATRAMPALLI-635 854**

**ANNA UNIVERSITY : CHENNAI 600 025**

**2022-2023**

1. **INTRODUCTION**
  - 1.1 Project Overview
  - 1.2 Purpose
2. **LITERATURE SURVEY**
  - 2.1 Existing Solution
  - 2.2 References
  - 2.3 Problem Statement Definition
3. **IDEATION & PROPOSED SOLUTION**
  - 3.1 Empathy Map Canvas
  - 3.2 Ideation & Brainstorming
  - 3.3 Proposed Solution
  - 3.4 Problem Solution fit
4. **REQUIREMENT ANALYSIS**
  - 4.1 Functional requirement
  - 4.2 Non-Functional requirements
5. **PROJECT DESIGN**
  - 5.1 Data Flow Diagrams
  - 5.2 Solution & Technical Architecture
  - 5.3 User Stories
6. **PROJECT PLANNING & SCHEDULING**
  - 6.1 Sprint Planning & Estimation
  - 6.2 Sprint Delivery Schedule
7. **WORKING WITH THE DATASET & DATA VISUALIZATION**
  - 7.1 Understanding the Dataset
  - 7.2 Loading the Dataset
  - 7.3 Visualization Chart
8. **CREATING THE DASHBOARD**
9. **ADVANTAGES & DISADVANTAGES**
10. **CONCLUSION**
11. **FUTURE SCOPZ**
12. **SOURCE CODE**
13. **GITHUB LINK**

# 1. INTRODUCTION

## 1.1 Project Overview

Bike share programs have risen in popularity in recent years and have been promoted as a lower carbon alternative to other forms of transit. Interest in bicycle sharing has been growing exponentially over the past decade, resulting in a proliferation of bike share systems in 712 cities across the world, encompassing 806,000 bicycles and 37,500 stations. This can be largely attributed to the successful incorporation of information technology in docking stations and mobile devices as well as improved logistics such as bicycle rebalancing to ensure responsive supply management. Cities often hope bike sharing will bring many benefits such as extending the reach of transit, substituting motorized trips, and encouraging non-cyclists to try cycling.

The premise of bicycle sharing is that it is a short-term bike rental system, based on varying timed memberships. Members of the bike share network have access to stations, consisting of a pay-station and multiple bike docks, across the system where bikes can be checked out from one station and returned to another nearest to their destination. The appeal of membership is 24/7 access to an automated bike rental network and utility of bikes in completing “last-kilometer connections” without the worry of storage or maintenance. The price system is set to encourage shorter trips (less than 30 minutes in time), with additional fees for any time used over that maximum.

There is evidence that bike share users switch to bike share from motorized transport, such as bus and auto, creating the potential for significant reductions in transportation related greenhouse gas or CO<sub>2</sub>e emissions. However, there is significant heterogeneity between different cities, showing that there is not a guaranteed CO<sub>2</sub>e reduction benefit from instituting bike share, especially if the trips would not have been made otherwise or are substituting walking and private bicycle trips.

## 1.2 Purpose

The purpose of this analysis is to create an operating report of Citi Bike for the year 2018. From this analysis, the following data visualizations will be created.

- 1.Total Number of Trips
- 2.What is Customer and subscriber with gender
- 3.Find the top bike used with respect to trip duration?
- 4.Calculating the number of bikes used by respective age groups.
- 5.Top 10 Start Station Names with respect to Customer age group

## 2. LITERATURE SURVEY

### 2.1 Existing Problem

**Spinlister** -Spinlister is an online hub for renting bikes from individuals or bike rental shops.

**Zagster** - Life is better on a bike! They are bringing bike share to communities across the USA.

**Motivate International** - Motivate is a global full-service bike share operator and technology innovator.

**Spin** - Spin is a stationless bike and electric scooter sharing service.

### 2.2 References

<https://craft.co/citi-bike/competitors>

Ines et al.,ScienceDirect-Social and Behavioral Sciences 111 ( 2014 ) 518 – 527  
“ Bicycle sharing systems demand”

Elias et al.,ScienceDirect Journal of Transport Geography 91 (2021)  
102971”What do trip data reveal about bike-sharing system users? “

FRANCESCO et al.,IEEE Access 2020”Bike Sharing and Urban Mobility in a  
Post-Pandemic”

“A long-term perspective on the COVID-19: The bike sharing  
system resilience under the epidemic environment”Journal of Transport &  
Health ,2021

Nguyen ThiHoai Thu, Chu Thi Phuong Dung, Vietnam 2017 International  
Conference on Advanced Technologies for Communications - Multi-source Data  
Analysis for Bike Sharing Systems

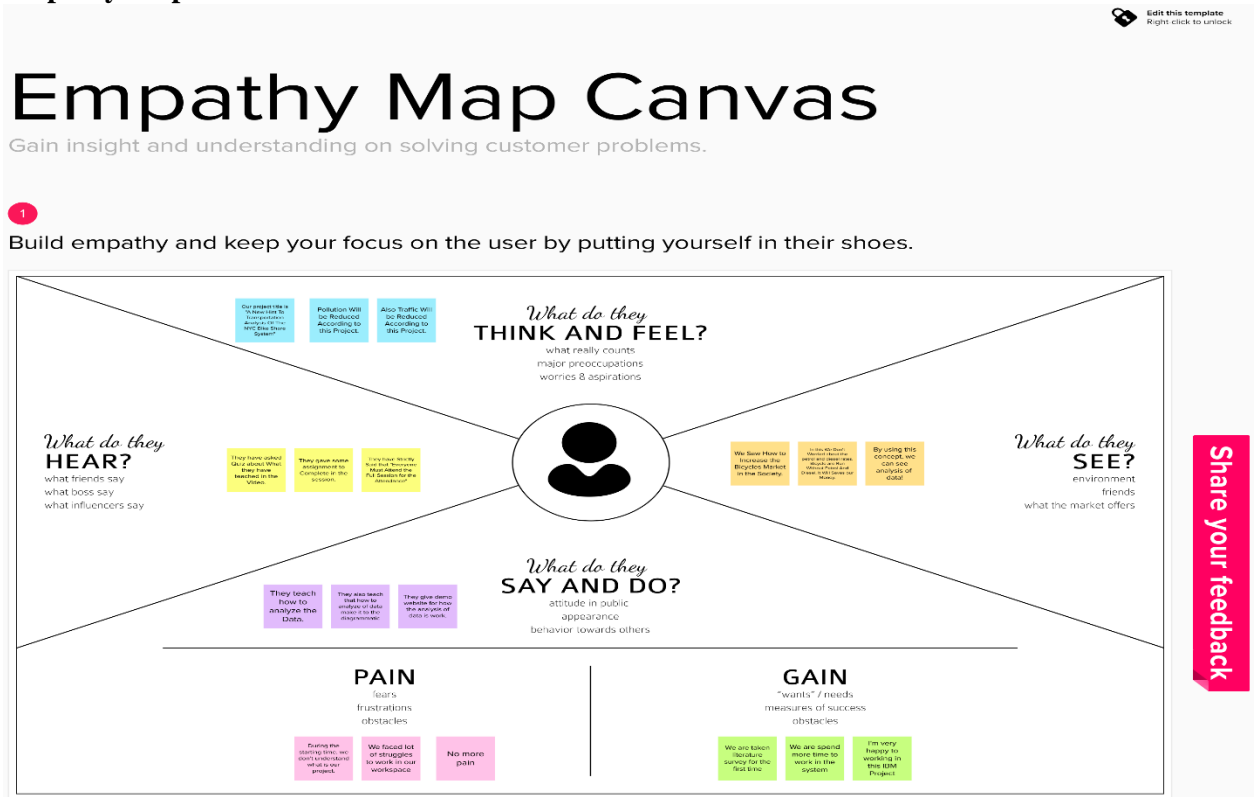
### 2.3 Problem statement Definition

In busy cities like New York the people are facing difficulties in analyzing the demand for bikes during peak hours.

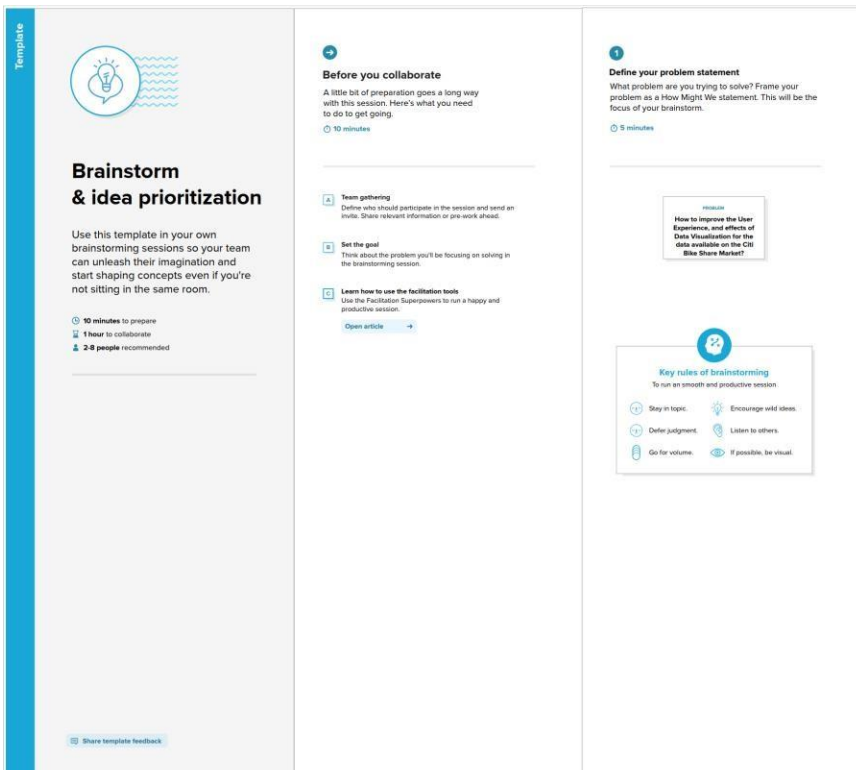
The main objective of this project is to predict bike patterns that will be extremely helpful for people to plan their travel.

### 3. IDEATION & PROPOSED SOLUTION

#### 3.1 Empathy Map Canvas



#### 3.2 Ideation and Brainstorming



2

## Brainstorm

Write down any ideas that come to mind that address your problem statement.

🕒 10 minutes

### TIP

You can select a sticky note and hit the pencil (switch to sketch) icon to start drawing!

**Thirupathi K**

Analysing the features in the dataset

**Ajay A**

Age-wise usage calculation of number of bikes

**Madhan R**

IBM Cognos Analytics Platform

**Vignesh P**

Three-tier architecture of User, Cognos and the Data

3

## Group ideas

Take turns sharing your ideas while clustering similar or related notes as you go. Once all sticky notes have been grouped, give each cluster a sentence-like label. If a cluster is bigger than six sticky notes, try and see if you can break it up into smaller sub-groups.

🕒 20 minutes

### TIP

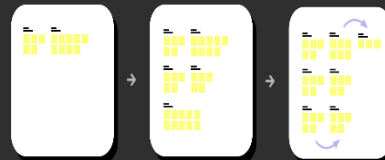
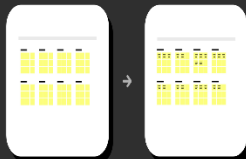
Add customizable tags to sticky notes to make it easier to find, to make groups, and to categorize important ideas as themes to tie your mind.

Preparation of the dataset and assignment to be done

Three-tier architecture of User, Cognos and the Data

Explore plots based on the bike usage count

Radial and Spiral Visualizations

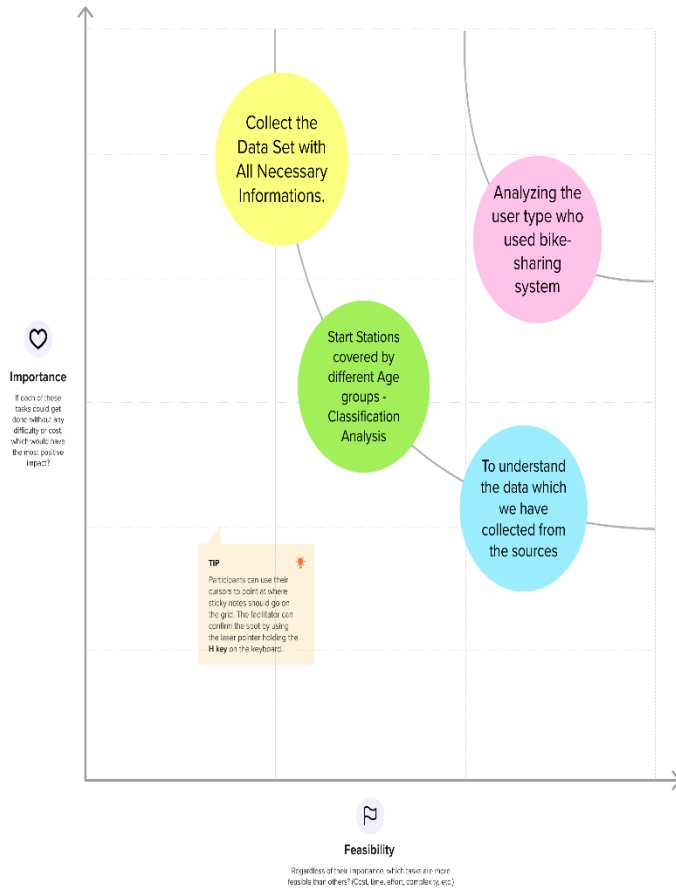


4

## Prioritize

Your team should all be on the same page about what's important moving forward. Place your ideas on this grid to determine which ideas are important and which are feasible.

20 minutes



→

## After you collaborate

You can export the mural as an image or pdf to share with members of your company who might find it helpful.

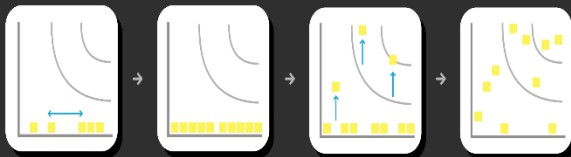
### Quick add-ons

- A Share the mural**  
Share a view link to the mural with stakeholders to keep them in the loop about the outcomes of the session.
- B Export the mural**  
Export a copy of the mural as a PNG or PDF to attach to emails, include in slides, or save in your drive.

### Keep moving forward

- Strategy blueprint**  
Define the components of a new idea or strategy.  
[Open the template →](#)
- Customer experience journey map**  
Understand customer needs, motivations, and obstacles for an experience.  
[Open the template →](#)
- Strengths, weaknesses, opportunities & threats**  
Identify strengths, weaknesses, opportunities, and threats (SWOT) to develop a plan.  
[Open the template →](#)

[Share template feedback](#)



### 3.3 Proposed Solution

S.No.	Parameter	Description
1.	Problem Statement (Problem to be solved)	The government needs a way to analyze the NYC bike share system so that they can enhance the system and give residents and visitors a fun, safe, affordable and convenient alternative to walking, taxis, buses etc.
2.	Idea / Solution description	The goal of this analysis is to create an operating report of Citi Bike for the year 2018. We are going to create different types of data visualizations using the various features of IBM Cognos Analytics so that the user can better understand the results of the analysis. It integrates reporting, modeling, analysis, dashboards etc so that the users can understand the available data, and make effective decisions. It includes predictive, descriptive, and exploratory techniques and provides an intuitive and straightforward interface that is easy to understand. Python's analytical functions can also be used for generating descriptive statistics and visualizations can also be created using Python's visualization libraries.
3.	Novelty / Uniqueness	Our solution gives faster results, reduces maintenance due to complete report coverage, and improved decision making - our reports and dashboards present the data in easily-understood formats.
4.	Social Impact / Customer Satisfaction	Bike share engages riders in physical activity, beneficial to health. In addition, it promotes green mobility and contributes to carbon neutrality. This analysis will help in understanding the association between bike share usage and the environment which is essential for system management and urban transportation planning.
5.	Business Model (Revenue Model)	This analysis might show that bike share is a relatively inexpensive and quick-to-implement urban transportation option compared to other transportation modes. The relative cost of launching a bikeshare system is less than investments in other transportation infrastructure, such as public transit and highways.
6.	Scalability of the Solution	This analysis presents evidence of the possible contribution of bike sharing systems to a more resilient transport system, as it can quickly provide alternative transport options to urban residents. As more data becomes available, particularly in other areas with identically comprehensive bike sharing systems, a clearer picture of the role of this transport mode in these emergency situations can be better evaluated by this analysis and provide results with an increased accuracy.



### 3.4 Problem Solution Fit

<p>Define CS, fit into</p> <p>Focus on J&amp;P, tap into BE, understand</p> <p>Identify strong TR &amp; EM</p>	<b>1. CUSTOMER SEGMENT(S)</b> <span>CS</span> Who is your customer?  Anyone who requires a cheap and efficient medium of transport for a short period of time with no need of maintenance.	<b>6. CUSTOMER</b> <span>CC</span> What constraints prevent your customers from taking action or limit their choices of solutions? i.e. spending power, budget, no cash, network connection, available devices.  A well maintained database with a clear info about the user and the bike and an availability of the steady internet connection should be ensured	<b>5. AVAILABLE SOLUTIONS</b> <span>AS</span> Which solutions are available to the customers when they face the problem or need to get the job done? What have they tried in the past? What pros & cons do those solutions have? i.e. pen and paper is an alternative to digital notetaking  Traditional way of manually documentable database could be maintained and shared. But there might be a possibility of human errors and confusion due to huge records.	Explore AS
	<b>2. JOBS-TO-BE-DONE / PROBLEMS</b> <span>J&amp;P</span> Which jobs-to-be-done (or problems) do you address for your customer? There could be more than one, explore different sides.  To get a detailed information and stats about the rented bike to ascertain a proper bike sharing system.	<b>9. PROBLEM ROOT CAUSE</b> <span>RC</span> What is the real reason that this problem exists? What is the back story behind the need to do this job? i.e. customers have to do it because of the change in regulations.  Manual accounting and tracking of the rented vehicles could sometimes result in loss in track of records of the current user and the bike, which may lead to some serious consequences that should be faced by the bike sharing service provider.	<b>7. BEHAVIOUR</b> <span>BE</span> What does your customer do to address the problem and get the job done? i.e. directly related: find the right solar panel installer, calculate usage and benefits; indirectly associated: customers spend too time on volunteering work (i.e. Greenpeace)  User help and support could be provided by including the customer care services in the interface and instruction manuals could also be provided to the each user of the rented bike to cross check and verify the working of the software, interface and the bike sharing system.	Focus on J&P, tap into BE, understand
	<b>3. TRIGGERS</b> <span>TR</span> What triggers customers to act? i.e. seeing their neighbour installing solar panels, reading about a more efficient solution in the news.  Common and a more advanced practices that encourages public and a shared transport medium.	<b>10. YOUR SOLUTION</b> <span>SL</span> If you are working on an existing business, write down your current solution first, fill in the canvas, and check how much it fits reality. If you are working on a new business proposition, then keep it blank until you fill in the canvas and come up with a solution that fits within customer limitations, solves a problem and matches customer behaviour.  The Solution will have an advanced tech improvements in the bike sharing system which would bring advancements in the society and will also could act as a factor that cuts CO2 emission	<b>8. CHANNELS OF BEHAVIOUR</b> <span>CH</span> <b>8.1 ONLINE</b> What kind of actions do customers take online? Extract online channels from #7  Steady network and an efficient database system should be made ensured.  <b>8.2 OFFLINE</b> What kind of actions do customers take offline? Extract offline channels from #7 and use them for customer development.  Ensure the proper working of bikes and the genuineness of the users.	Extract online & offline CH of BE

## 4. REQUIREMENT ANALYSIS

### 4.1 Functional Requirement

FR No.	Functional Requirement (Epic) Sub Requirement (Story / Sub-Task)
FR-1	Collection of user data Lyft citi bike's official website provides the data to help with analysis, development, visualization etc. Data is collected from these published files.
FR-2	Analysing the user data This data is used as input for creating various types of visualizations and analysis is done and a dashboard is created.
FR-3	Display the data The dashboard is used to display the top bike used with respect to trip duration, top 10 Start Station Names with respect to customer age group, to find the customer and subscriber with gender, to find total number of trips & calculating the number of bikes used by respective age groups.

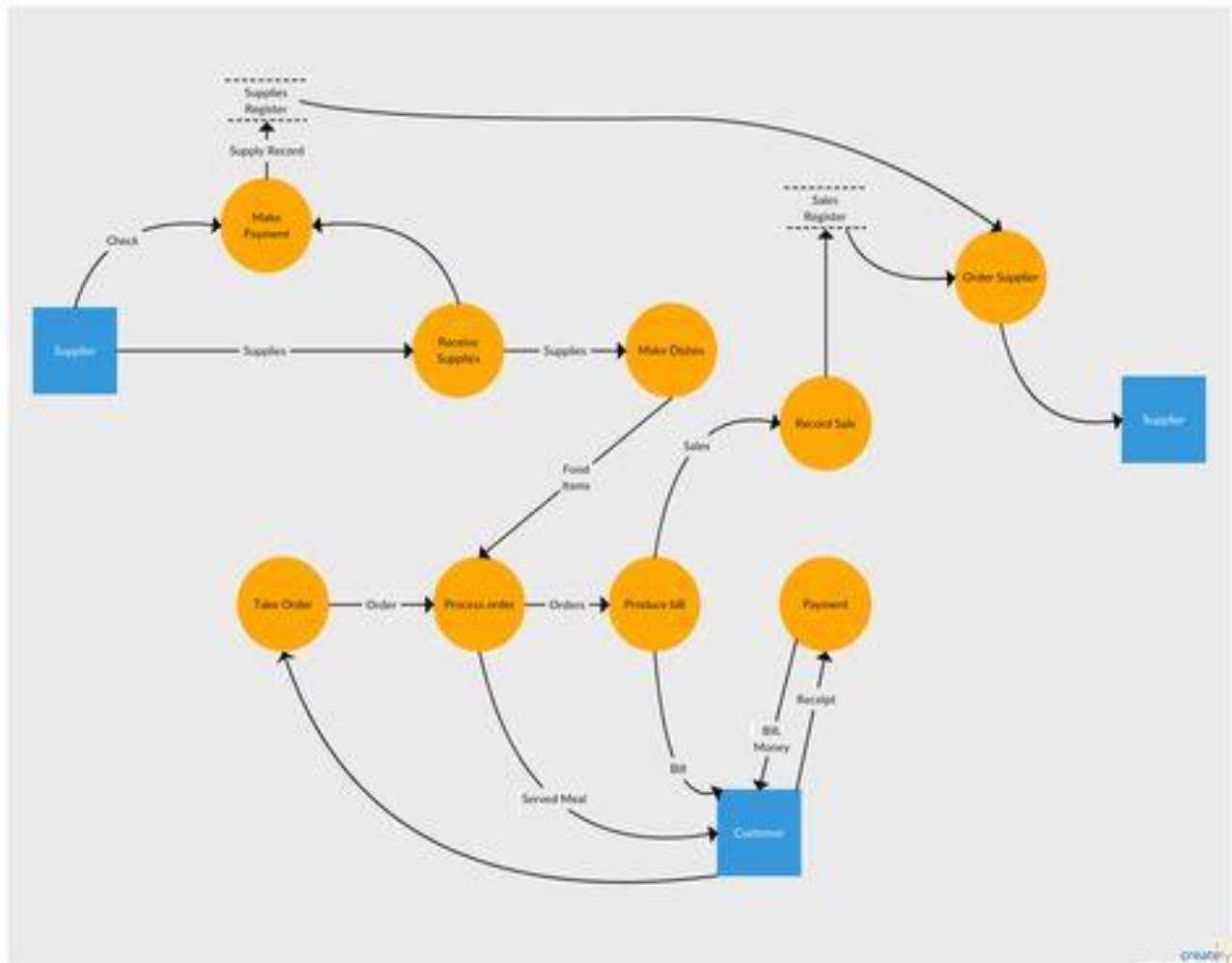
### 4.2 Non-Functional Requirement

FR No.	Non-Functional Requirement	Description
NFR-1	Usability	This dashboard provides an easily understandable report which facilitates many people and tourists who use bicycles to complete their work and enjoy themselves. It provides many benefits such as measures data like distance, and help with tasks such as route planning, expansion of the bicycle sharing system, manufacturing of desired bikes etc. The benefits of Bicycle sharing systems could be reduced vehicle emissions, reduces energy consumption, improve health benefits, financial savings for individuals, reduced congestion and

		fuel consumption.
<p><b>NFR-2 Security</b> The citi bike usage data is secured with appropriate caution as crucial decisions will be made based on this data. We can restrict access to this data and the visualization reports.</p>		
NFR-3	<b>Reliability</b>	This analysis provides a reliable and an efficient way to grasp on the performance of the citi bike sharing system in the year 2018. It makes use of the available dataset precisely and gives accurate data visualizations that can be used to improve the citi bike sharing system.
NFR-4	<b>Performance</b>	Performance of bike sharing system is defined as operational efficiency and spatial effectiveness of bike sharing system. The operational efficiency of bike sharing system aims at understanding the characteristics of public bike users, and evaluating the conditions of bike lanes from the perspective of public bike users. The effectiveness of bike sharing system dashboard aims at analyzing the characteristics of bike stations, and accessibility between bike stations and other facilities. The evaluation results can be used to improve the public bicycle sharing program.
NFR-5	<b>Availability of bikes</b>	A bicycle-sharing system is a shared transport service where bicycles are available for shared use by individuals for a short-term at low or zero cost. The programs themselves include both docking and dockless systems, where docking systems allow users to borrow a bike from a dock and return at another node or dock within the system — and dockless systems, which offer a node-free system relying on smart technology. In either format, systems may incorporate smartphone web mapping to locate available bikes and docks.
NFR-6	<b>Scalability</b>	This analysis presents evidence of the possible contribution of bike sharing systems to a more resilient transport system, as it can quickly provide alternative transport options to urban residents. As more data becomes available, particularly in other areas with identically comprehensive bike sharing systems, a clearer picture of the role of this transport mode in these emergency situations can be better evaluated by this analysis and provide results with an increased accuracy.

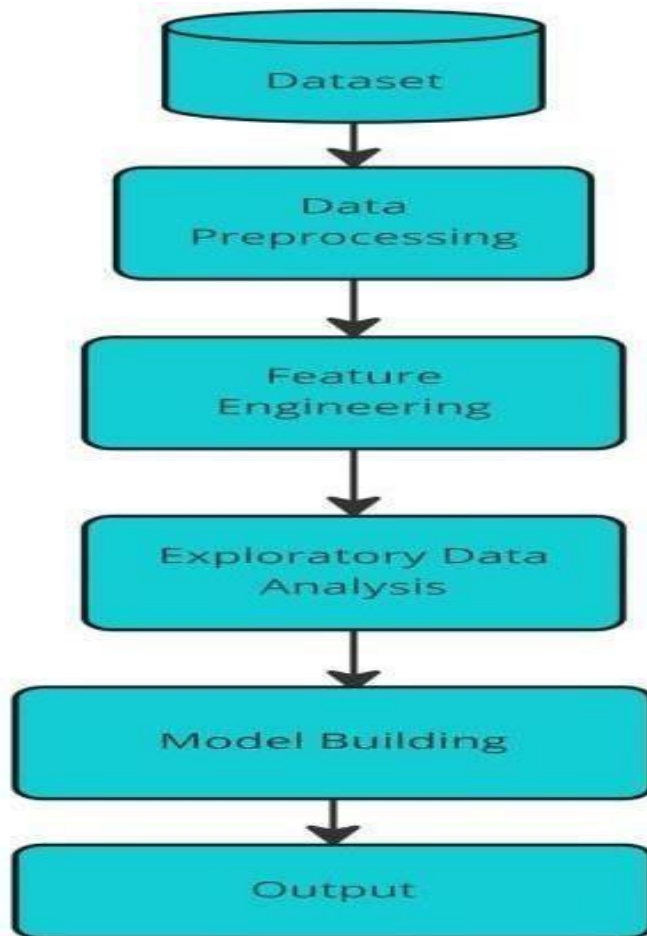
## 5. PROJECT DESIGN

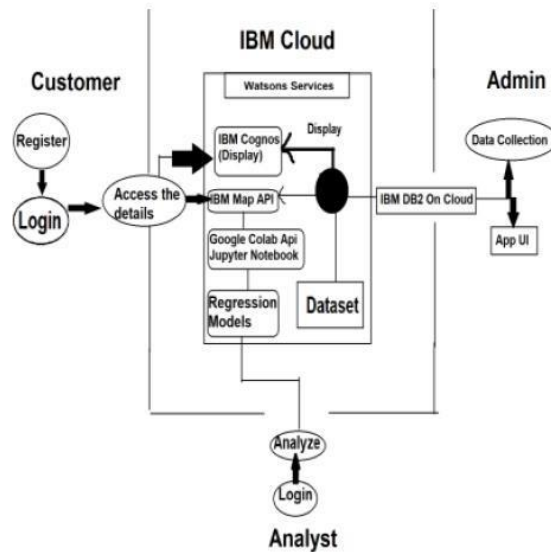
### 5.1 Data Flow Diagram



User Type	Functional Requirement (Epic)	User Story Number	User Story / Task	Acceptance criteria	Priority	Release
Customer(Analysts at Citi, Government)	Registration	USN-1	As a user, I should be able to register to see the dashboard as a new user	Successful Registration	High	Sprint-1
Customer(Analysts at Citi, Government)	Login	USN-2	As a user I should be able to login to see the dashboard with the correct credentials	Successful Login with correct credentials	High	Sprint-1
Customer(Analysts at Citi, Government)	Accessing the dashboard	USN-3	As a user, I should be able to view the visualizations displayed	Should be able to view the following analysis among others : <ol style="list-style-type: none"> <li>1. Total number of trips</li> <li>2. Subscriber and Customer with gender</li> <li>3. Top Bike used with respect to duration</li> <li>4. Number of bikes used by different age groups</li> <li>5. Top start station name with respect to customer age group</li> </ol>	High	Sprint-1
Customer(Analysts at Citi, Government)	Manipulating the data	USN-4	As a user I should be able to apply some modifications to the data to see how the resultant visualizations change	I should have the permission to manipulate the data	High	Sprint-2

## 5.2 Solution & Technical Architecture





**Table-1 :Components & Technologies:**

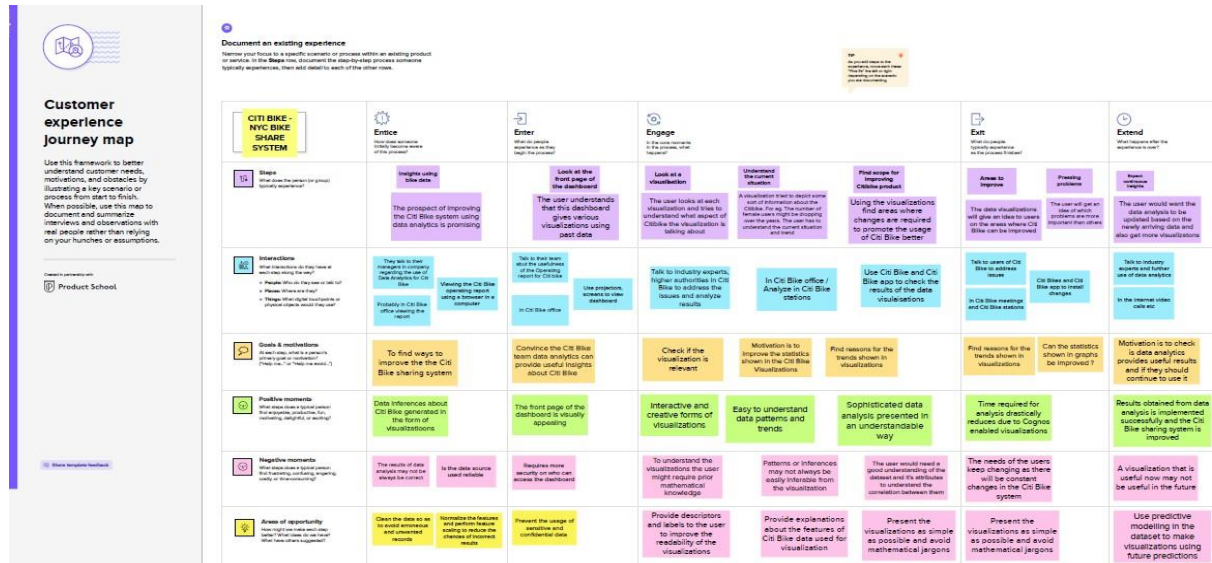
S.No	Component	Description	Technology
1.	User Interface	User can Interact with web Application	HTML, CSS, JavaScript .
2.	Data Pre processing	Pre processing of data should be done	Python
3.	Feature Engineering	Feature engineering of Dataset by adding new values to the existing dataset.	Python
4.	Exploratory Data Analysis	Exploring the data using boxplot, pie plot, scatter plot etc..	Python
5.	Model Building	Build the model using machine learning algorithm	python
6.	Data Storage	Database Service on Cloud	IBM DB2, IBM Cloudant etc.
7.	User Interface	Dashboard showing the details of the trip duration, no of trips, bike usage etc..	HTML,CSS, JavaScript.

**Table-2: Application Characteristics:**

S.No	Characteristics	Description	Technology
1.	Security Implementations	The main security concern is for users account hence proper login mechanism should be used to avoid hacking.	e.g. SHA-256, Encryptions, IAM Controls, OWASP etc.
2.	Availability	The system shall be available 24 hours a day 7 days a week. User can access at anytime	
3.	Performance	The system should require a fair amount of speed especially while browsing through the catalogue	

### 5.3 User Stories





## 6. PROJECT PLANNING & SCHEDULING

### 6.1 Sprint Planning & Estimation

Sprint	Functional Requirement (Epic)	User Story Number	User Story / Task	Story Points	Priority	Team Members
Sprint-1	Data Preparation	USN-1	As an analyst, find and extract the Citi-bike data for the year 2018 from the given bucket list of datasets.	3	Medium	Mahesh S M
Sprint-1	Data Preprocessing	USN-2	As an analyst, upload the extracted dataset to IBM Cognos platform.	1	Medium	Pranav M
Sprint-1	Data Preprocessing	USN-3	As an analyst, understand the dataset working with and give a brief overview of what each feature represents.	1	Low	Kabilraj S
Sprint-2	Data Preprocessing	USN-4	As an analyst, prepare the data for analysis by handling missing values.	2	Medium	Robert Samuel J
Sprint-2	Analysis	USN-5	As an analyst, perform Exploratory Data Analysis on the filtered dataset to identify patterns and relationships between various features present.	3	High	Mahesh S M, Pranav M
Sprint-2	Visualization	USN-6	As an analyst, create 5 various visualizations charts using IBM Cognos	10	High	Kabilraj S, Robert Samuel J
Sprint-3	Visualization	USN-7	As an analyst, creation of a dashboard by the created Visualizations charts to understand business insights making it a interactive dashboard.	10	High	Pranav M, Kabilraj S
Sprint-3	Visualization	USN-8	As an analyst, prepare report and story by applying predictive analytics to enhance visualizations chart interactions.	8	Low	Mahesh S M, Robert Samuel J
Sprint-3	Exporting	USN-9	As an analyst, export the analytics to share the work to showcase to others.	2	Medium	Mahesh S M, Kabilraj S
Sprint-4	Testing	USN-10	As an analyst to perform an user acceptance testing.	3	Medium	Mahesh S M, Pranav M
Sprint-4	Delivery	USN-11	As an analyst to prepare project report, demo video, ensuring all project deliverables are submitted and available.	5	High	Mahesh S M, Robert Samuel
Sprint-4	Registration	USN-12	As a user, should view the web application of the interactive dashboard.	2	Low	Pranav M, Robert Samuel J





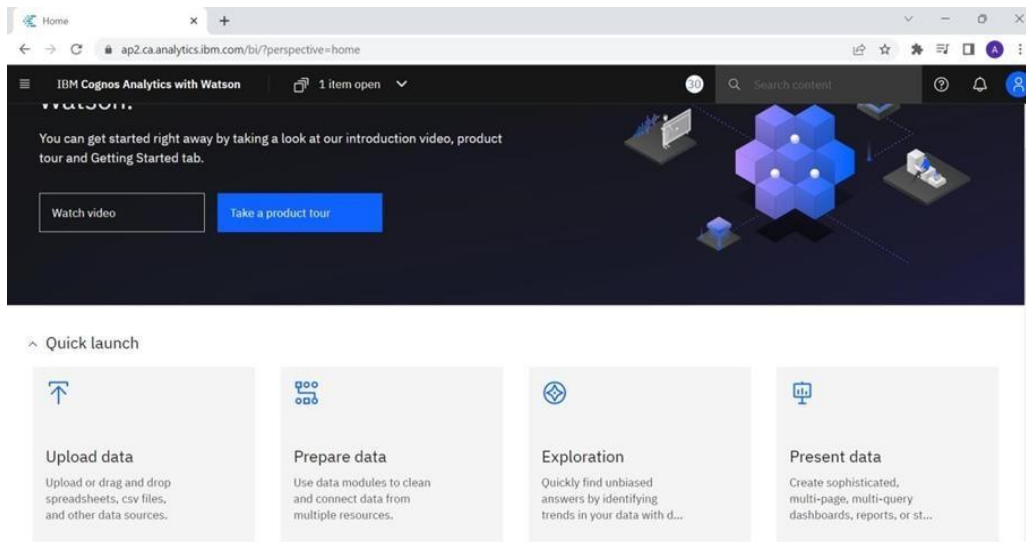
**Dataset Link:** [Dataset](#)

1. Trip Duration: How long a trip lasted in seconds
2. Start Date and Time: EX->01-06-2013 00:00:01
3. Stop Date and Time: EX->01-06-2013 00:11:36
4. Start Station ID: Unique identifier for each station
5. Start Station Name
6. Start Station Latitude: Coordinates
7. Start Station Longitude: Coordinates
8. End Station ID: Unique identifier for each station
9. End Station Name
10. End Station Latitude
11. End Station Longitude
12. Bike ID: Unique identifier for each bike
13. User Type (Customer = 24-hour pass or 3-day pass user; Subscriber = Annual Member):  
Customers are usually tourists, subscribers are usually NYC residents
14. Year of Birth: Self-entered, not validated by an ID

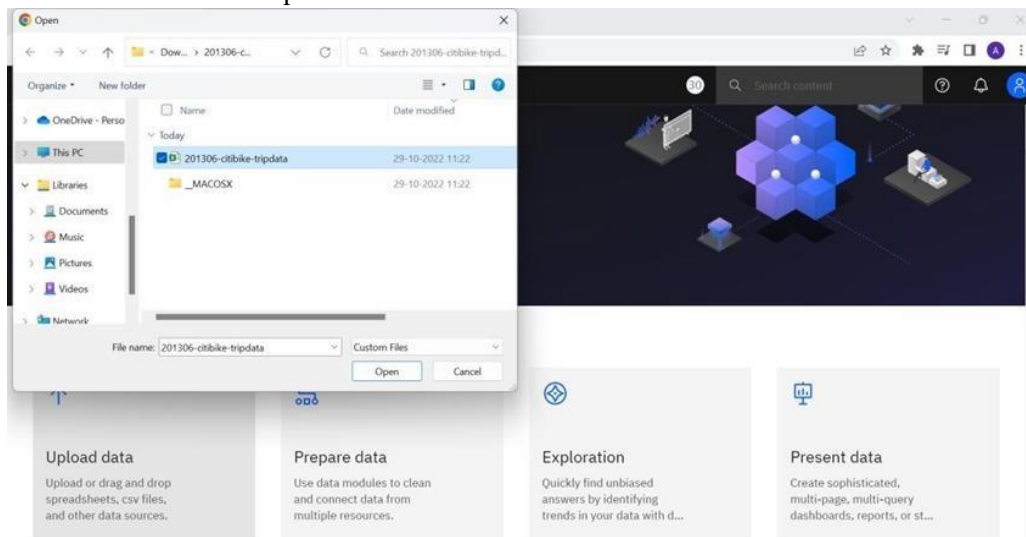
Gender (Zero=unknown; 1=male; 2=female): Usually unknown for customers since they often sign up at a kiosk

## **7.2 Loading the dataset**

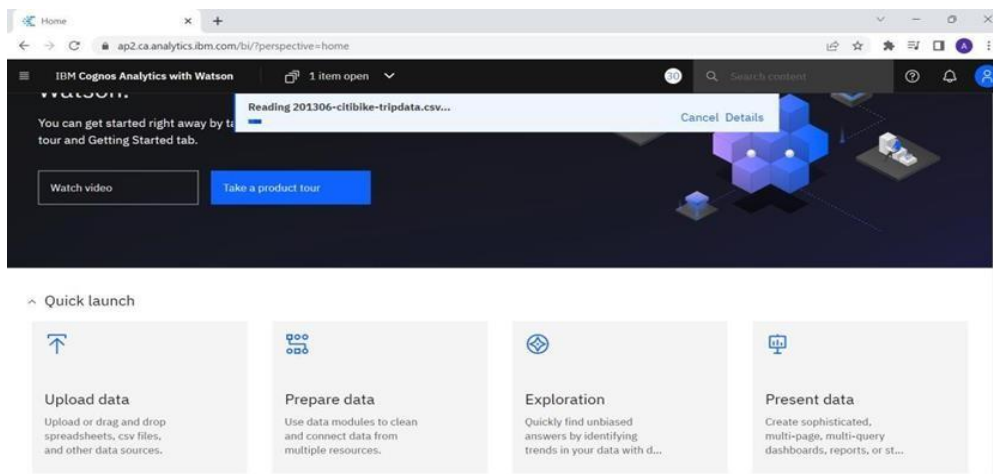
Open Cognos Analytics and click upload data



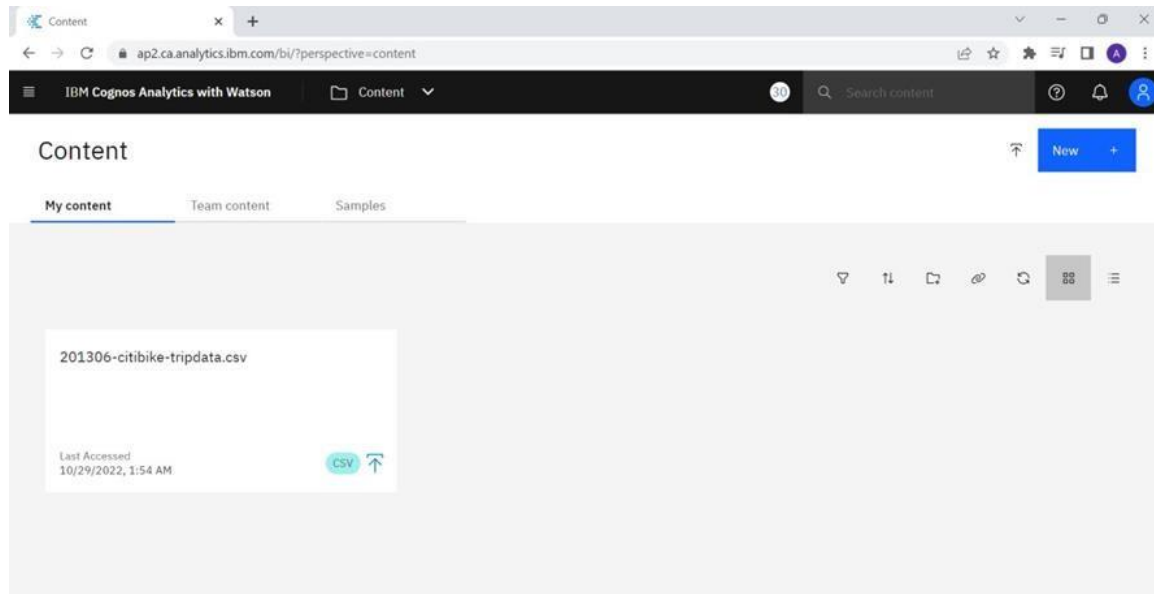
Select the dataset to be uploaded



The excel file is getting uploaded in Cognos Analytics

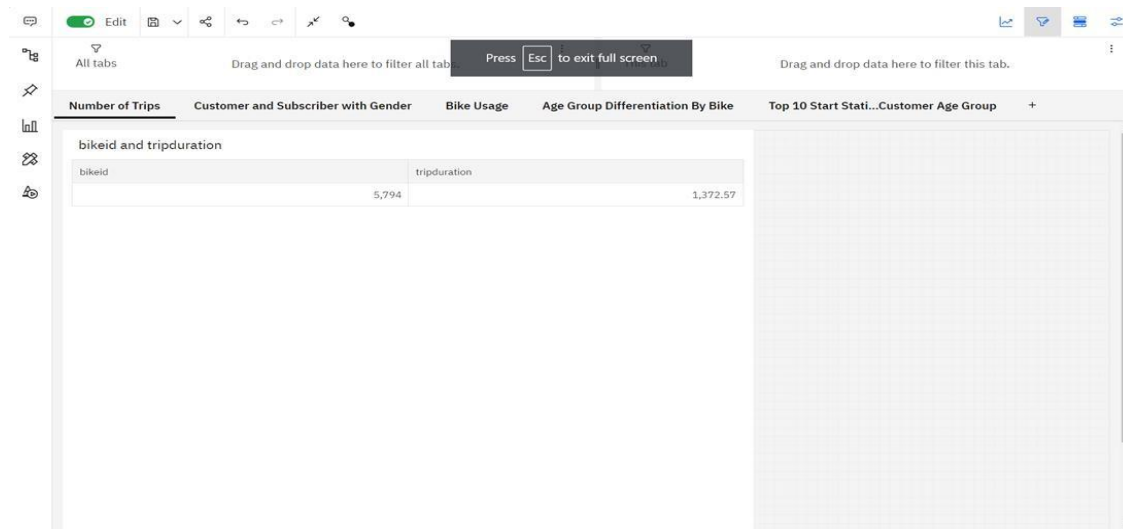


The dataset can be accessed in My Content in Cognos Analytics

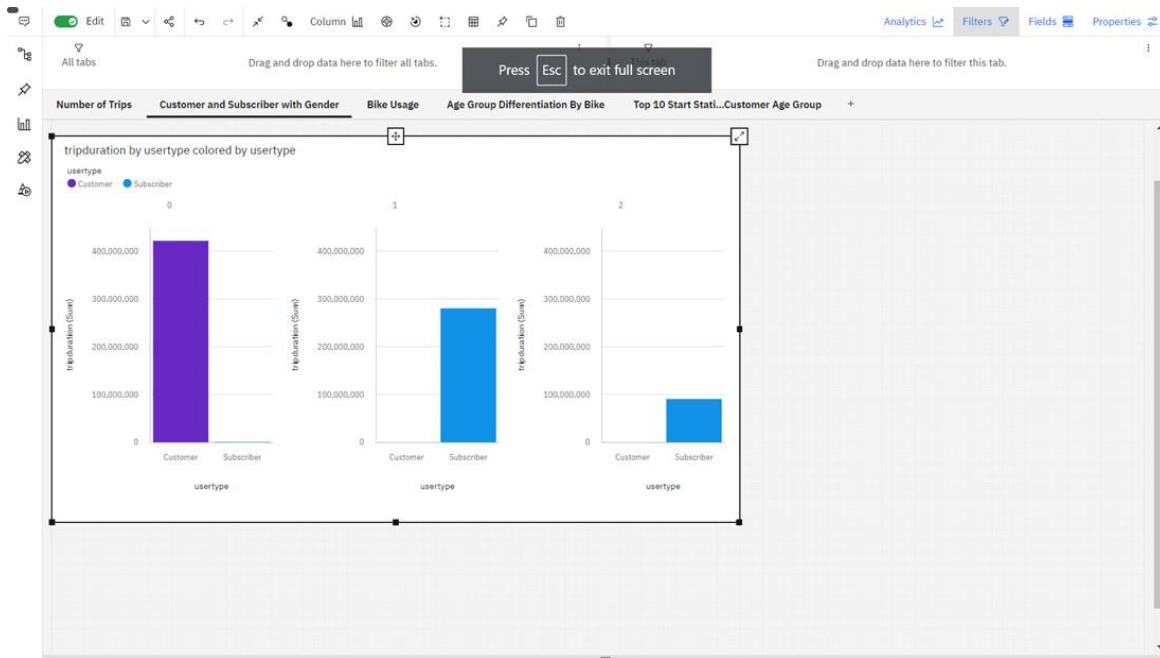


## 7.3 Visualization charts

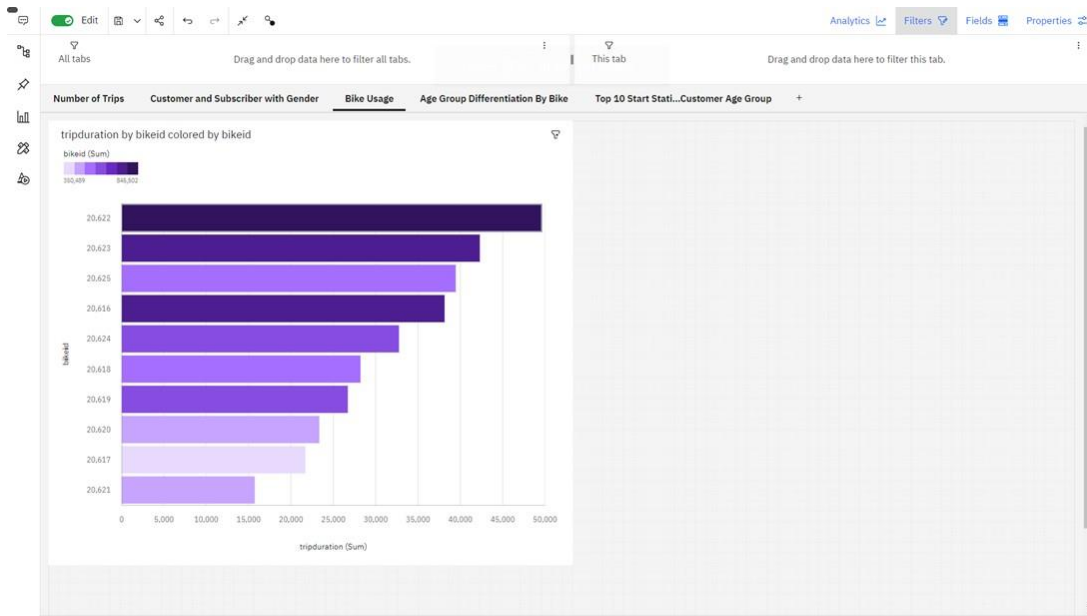
Number of Trips:



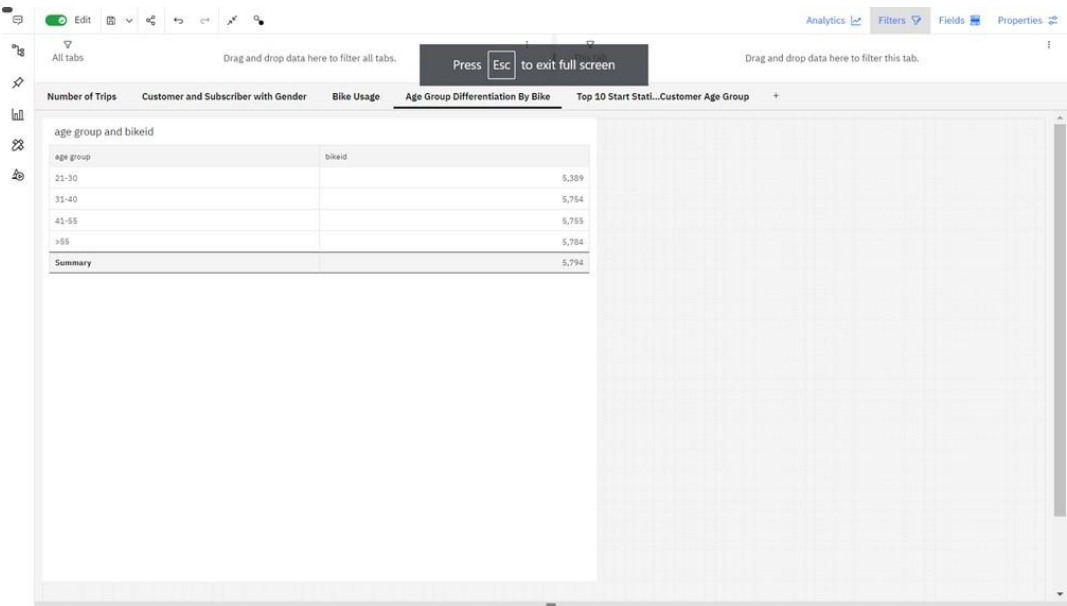
## Customer and Subscriber with Gender:



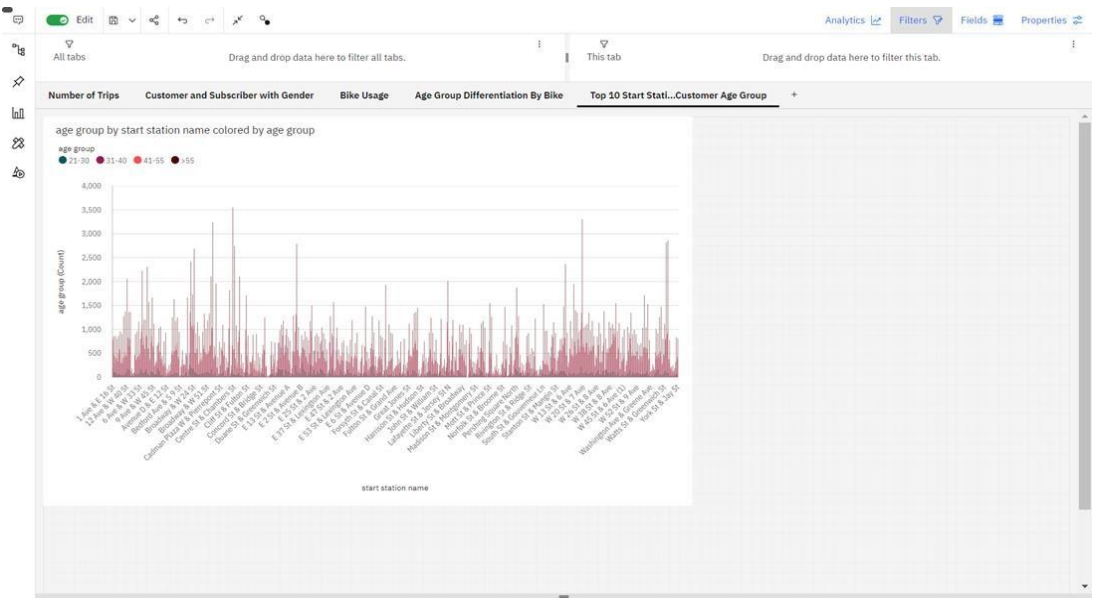
## Bike Usage:



Age group differentiation by bike:



Top 10 Start Station Names with Respect to Customer Age Group:



Gender Variation





## **6. ADVANTAGES AND DISADVANTAGES**

The benefits of bike sharing schemes include transport flexibility, reductions to vehicle emissions, health benefits, reduced congestion and fuel consumption, and financial savings for individuals.

One can easily analyze and understand trends in bike sharing patterns with the created dashboard. With no prior skills and knowledge about the tools that we use for analysis, anyone (literate or illiterate) can easily infer the knowledge that we represent in various charts or graphs or maps. So that it would be helpful to users and companies to make appropriate decisions in the future.

## **7. CONCLUSION**

Based on the quantitative as well as visual analysis of the New York bike share system, a number of interesting insights were gained.

One obvious conclusion was that there is a strong seasonal variation in the system usage with maximum usage in summer and minimum usage in winter. This was initially hypothesized because of the harshness of New York's harsh winters and the treacherous riding conditions that exist during that time. However, despite the adverse weather conditions, there is a strong core demographic that consistently uses the system. This conclusion is based on that fact that even during the months of January and February which are the peak winter months, there are more than two hundred thousand trips in the system

New York has a strong public transit system, and the bike share system seems to complement it quite well with a majority of the highest used stations located either close to subway lines or the commuter rail stations in the city

Based on the locations of the stations and the duration of trips, it can be hypothesized that bike shares are replacing last mile trips that would otherwise be done either on foot or on public transit. This is particularly true in case of New York where a combination of dense public transit network, the road congestion during peak hours and the average trip distance as calculated create a situation where the only potential trips that the bike share system is replacing currently are those that would otherwise have been undertaken either on foot or on public bus.

## 8. FUTURE SCOPE

NYC is a very crowded and happening place which leads to lots of pollution. And in this busy world people are always worried about transportation this bike sharing system reduces that stress. With increase in population pollution also increases. So it is in our hands to reduce pollution and to make a better future for our younger generations. We can analyze which station needs more bikes and any area needs new station to be installed. The survey outcomes indicates the needs for improved techniques in bike sharing analytics. There exists a lot of scope in this research area.

## 9. SOURCE CODE

```
#%% md

# SPRINT **3**

#%%

import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import plotly.express as px
from datetime import datetime
from pprint import pprint

from pydrive.auth import GoogleAuth
from pydrive.drive import GoogleDrive
from google.colab import auth
from oauth2client.client import GoogleCredentials

#%%

path = "/content/dataset.csv"
```



```
df = pd.read_csv(path)
print(df)
```

```
## % %
```

```
df.head()
```

```
## % %
```

```
df.describe()
```

```
## % %
```

```
df.info()
```

```
## % %
```

```
df.isnull().sum()
```

```
## % %
```

```
df[df['starttime'].isnull()]
```

```
## % %
```

```
df[df['stoptime'].isnull()]
```

```
## % %
```

```
df = df[:-1]
```

```
## % %
```

```
df.isnull().sum()
```

```
## % %
```

```
print(type(df["start station latitude"][0]))
print(df["start station latitude"][0])
```

```
## % %
```

```
df['start station name'].unique()
```

```
## % %
```

```
def camel_case(city):
    try:
        city = city.split(' ')

```

```

    city = ' '.join([x.lower().capitalize() for x in city])
    if city == 'Unknown':
        return np.nan
    else:
        return city
except:
    return np.nan

```

```

# Apply camel_case function to City column
df['start station name'] = df['start station name'].apply(camel_case)
df['start station name'].value_counts()

```

```

# %%

```

```

df.count()

```

```

# %%

```

```

df["tripduration"] = pd.to_numeric(df["tripduration"])
res = df.iloc[52323]
print(res["tripduration"])

```

```

# %%

```

```

df_filtered = df[df['tripduration'] != "tripduration"]
df_filtered["tripduration"] = pd.to_numeric(df_filtered["tripduration"])

```

```

df = df_filtered
type(df["tripduration"][0])

```

```

# %%

```

```

type(df["start station latitude"][0])

```

```

# %%

```

```

type(df["end station longitude"][0])

```

```

# %%

```

```

type(df["bikeid"][0])

```

```

# %%

```

```

type(df["birth year"][0])

```

```

# %%

```

```

type(df["gender"][0])

```

```

# %%

```

```
type(df["starttime"][0])
```

```
# % %
```

```
df["starttime"] = pd.to_datetime(df["starttime"])
```

```
df["stoptime"] = pd.to_datetime(df["stoptime"])
```

```
type(df["starttime"][0])
```

```
# % %
```

```
df["starttime"][0] < df["stoptime"][0]
```

```
# % %
```

```
df.info()
```

```
# % %
```

```
def find_outliers_IQR(df):
```

```
    q1=df.quantile(0.25)
```

```
    q3=df.quantile(0.75)
```

```
    IQR=q3-q1
```

```
    outliers = df[((df<(q1-1.5*IQR)) | (df>(q3+1.5*IQR)))]
```

```
    return outliers
```

```
outliers = find_outliers_IQR(df["birth year"])
```

```
print('number of outliers: ' + str(len(outliers)))
```

```
print('max outlier value: ' + str(outliers.max()))
```

```
print('min outlier value: ' + str(outliers.min()))
```

```
# % %
```

```
df["gender"].value_counts()
```

```
# % %
```

```
temp_df = df[df["birth year"] <= 1957]
```

```
temp_df["gender"].value_counts()
```

```
# % %
```

```
df.shape
```

```
# % %
```

```
df.to_csv('cleaned_dataset.csv', index=False)
```

```
# % % md
```

```
# **SPRINT 4**
```

```
# % %
```

```
path = "/content/cleaned_dataset.csv"
edadf = pd.read_csv(path)
print(edadf)
```

```
# %%
```

```
temp = edadf
```

```
# %%
```

```
temp.head()
```

```
# %%
```

```
temp.describe()
```

```
# %%
```

```
temp.info()
```

```
# %%
```

```
temp["starttime"] = pd.to_datetime(temp["starttime"])
temp["stoptime"] = pd.to_datetime(temp["stoptime"])
temp.info()
temp["Hour"] = temp["stoptime"].dt.hour - temp["starttime"].dt.hour
temp.head()
```

```
# %%
```

```
temp.shape
```

```
# %%
```

```
temp['Age'] = 2022 - temp['birth_year']
temp.head()
```

```
# %%
```

```
Age_Groups = ["<20", "20-29", "30-39", "40-49", "50-59", "60+"]
Age_Groups_Limits = [0, 20, 30, 40, 50, 60, np.inf]
Age_Min = 0
Age_Max = 100
temp["Age_group"] = pd.cut(temp["Age"], Age_Groups_Limits, labels=Age_Groups)
temp.head()
```

```
# %%
```

```
trips_df = pd.DataFrame()
```

```

trips_df = temp.groupby(['start station name','end station name']).size().reset_index(name =
'Number of Trips')
trips_df = trips_df.sort_values('Number of Trips',ascending = False)
trips_df['start station name'] = trips_df['start station name'].astype(str)
trips_df['end station name'] = trips_df['end station name'].astype(str)
trips_df['Routes'] = trips_df['start station name'] + " to " + trips_df['end station name']
trips_df = trips_df[:50]
trips_df = trips_df.reset_index()
trips_df

#%/%

px.pie(values = temp['gender'].value_counts(),
names =temp['gender'].value_counts().index,
title ="Gender Variation")

#%/%

px.bar(x=temp["start station name"].value_counts().index,
y=temp["start station name"].value_counts().values,
labels={'x':'Start Station Name','y':"Count"})

#%/%

px.bar(x=temp["end station name"].value_counts().index,
y=temp["end station name"].value_counts().values,
labels={'x':'End Station Name','y':"Count"})

#%/%

px.bar(x=temp["Hour"].value_counts().index,
y=temp["Hour"].value_counts().values,
title = "Hour usage of Citi Bikes",
labels={'x':'Time','y':"Number of people using bike"})

```

## 10. GITHUB LINK

<https://github.com/IBM-EPBL/IBM-Project-31013-1660194362>