

Project Report

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 SCOPE**
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 CO₂e emissions. However, there is significant heterogeneity between different cities, showing that there is not a guaranteed CO₂e 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

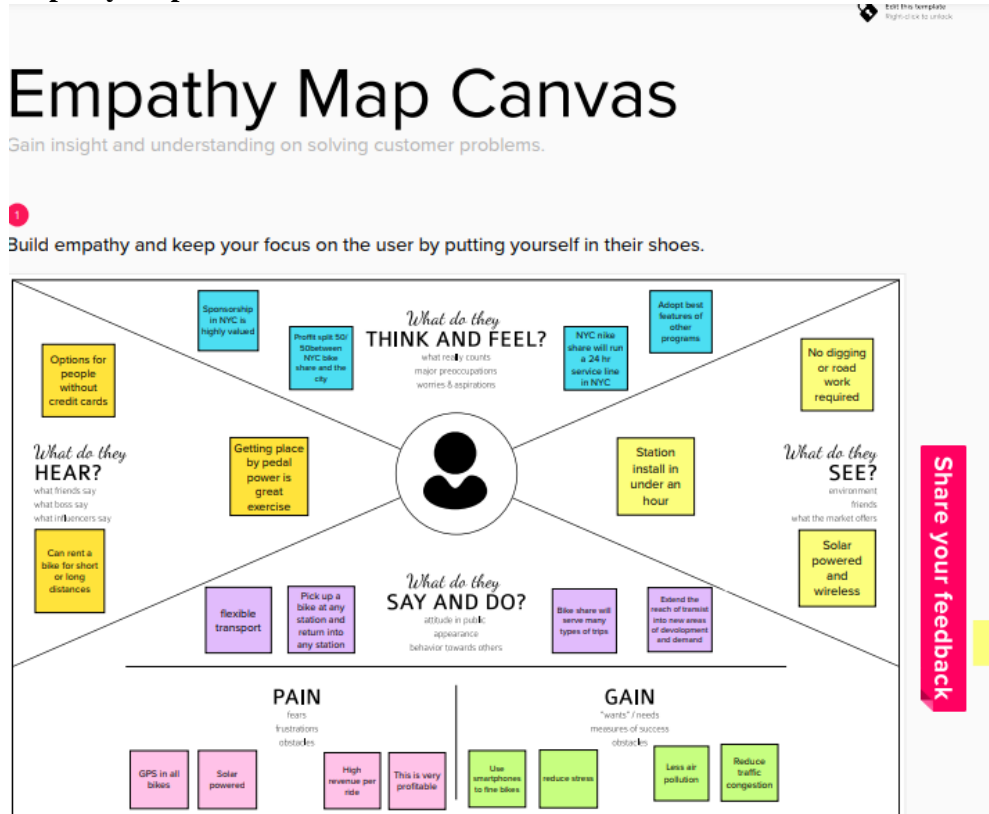
2.3 Problem statement Definition

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

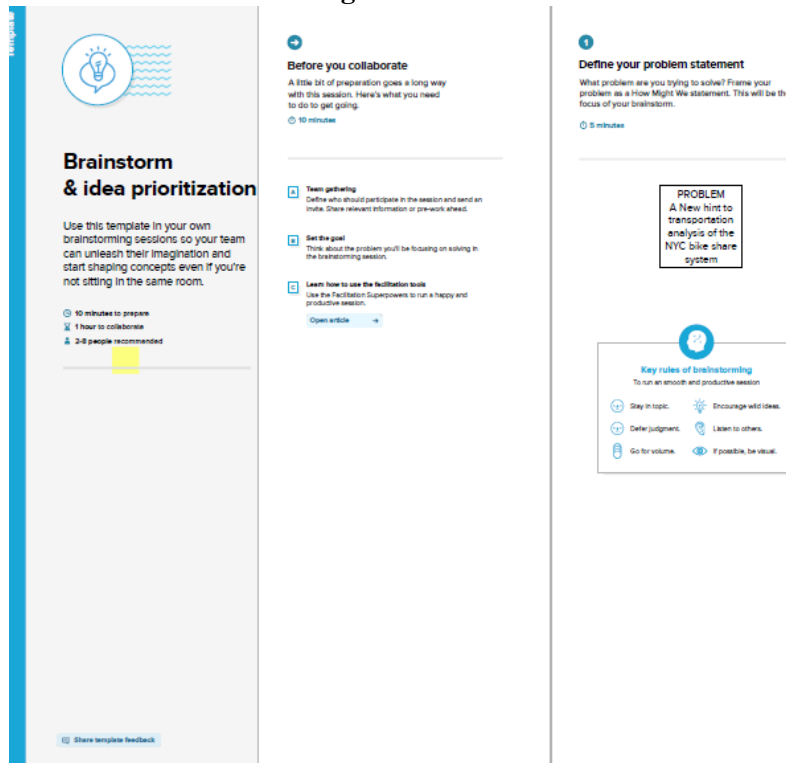
The main objective of this project is to predict bike patterns that will be extremely helpful for people to plan their travel.

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 button to start editing.

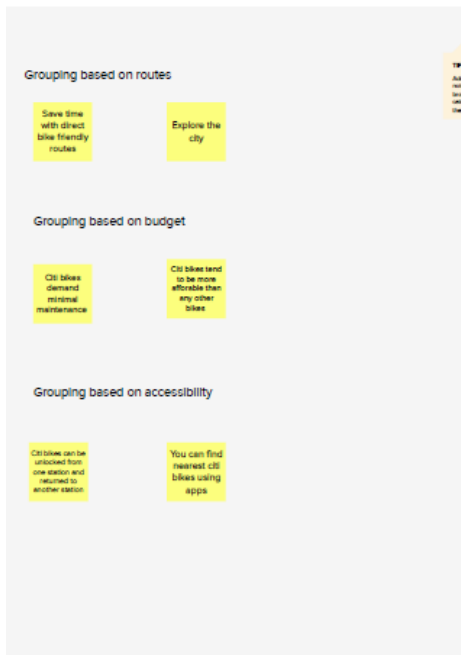


3

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



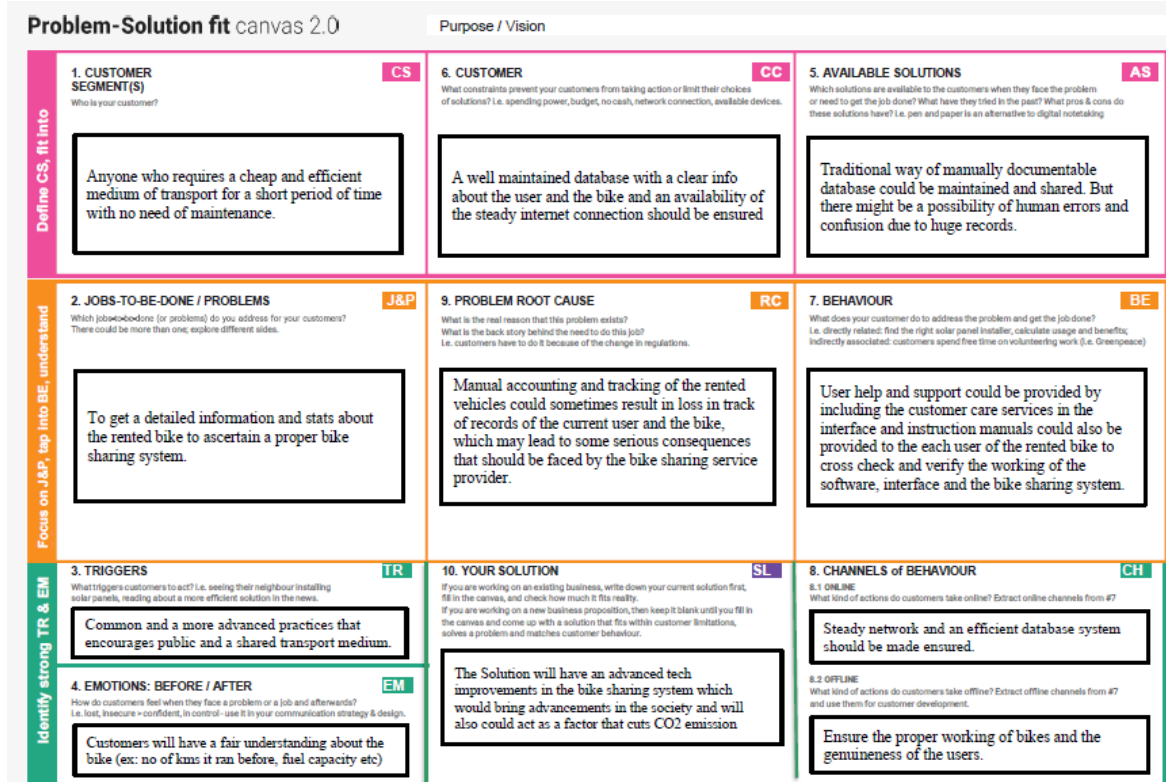
TIP
Add a name to each group using labels

3.3 Proposed Solution

S.NO	Parameter Description
1	<p>Problem Statement (Problem to be solved)</p> <ol style="list-style-type: none">1. Bike share programs have risen in popularity in recent years and have been promoted as a lower carbon alternative to other forms of transit.2. The premise of bicycle sharing is that it is a short-term bike rental system, based on varying timed memberships3. The trips would not have been made otherwise or are substituting walking and private bicycle trips
2	<p>Idea / Solution description</p> <ol style="list-style-type: none">1. The planning process for the Citi Bike program established an open door policy, encouraging input early and often from the citizens of New York City.2. Low-income people are less likely than middle- and upper-income people to have a credit card.3. Sites should ensure maximum visibility and access.4. Sites must not impede the use of any existing facilities, such as bus stops or fire hydrants.
3	<p>Novelty / Uniqueness</p> <ol style="list-style-type: none">1. Transport flexibility2. Reductions to vehicle emissions3. Health benefits4. Reduced congestion and fuel consumption5. Financial savings for individuals.

4	<p>Social Impact / Customer Satisfaction</p> <ol style="list-style-type: none"> 1. Transportation 2. Recreation of cycling 3. Enjoyable sport 4. Low cost
5	<p>Business Model (Revenue Model)</p> <ol style="list-style-type: none"> 1. Zero deaths since it launched. 2. Lack of public subsidies 3. Battery-powered bikes 4. The model is trained using open
6	<p>Scalability of the Solution</p> <ol style="list-style-type: none"> 1. Improved customer engagement 2. Reduce customer acquisition cost 3. Economical development 4. Immediate response for customer queries

1.1 Problem Solution Fit



2. REQUIREMENT ANALYSIS

2.1 Functional Requirement

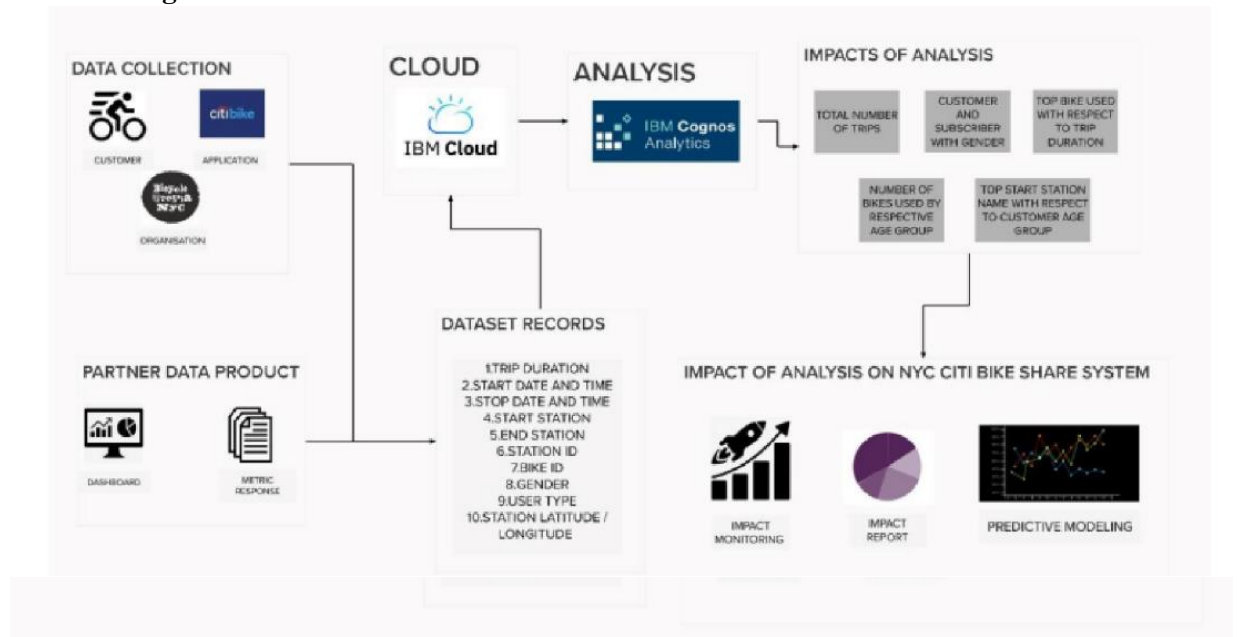
FR No.	Functional Requirement (Epic)	Sub Requirement (Story / Sub-Task)
FR-1	User Registration	Registration through Form Registration through Gmail Registration through LinkedIn
FR-2	User Confirmation	Confirmation via Email Confirmation via OTP
FR-3	Collection of Data	Utilizing the NYC Citi Bike assists in gathering information on the various trips that various users of Citi Bike take. These data were then organised into datasets and made available for further study and visualisation.
FR-4	Analysis of Data	Preprocessing and filtering the provided data in accordance with the sub-requirements task's is part of the analysis process. Data analysis and visualisation are both aided by the use of machine learning algorithms to glean more insights from the data..
FR-5	Display (Visualization) of Data	Various visualisations are used depending on the sub-task being handled. These visualisations are then combined and shown on a dashboard, which is a tool for giving customers business information. Finding the top 10 Start stations according to customer age group and showing the most popular bikes according to trip time are a few of the various sub-tasks included in this requirement.

2.2 Non-Functional Requirement

FR No.	Non-Functional Requirement	Description
NFR-1	Usability	The dashboard gives users access to an operational report that is simple to read and useful for understanding market trends and company insights. Data can be examined from various angles and in more depth by using an interactive dashboard to drill down and filter operating information.
NFR-2	Security	Based on the Citi Bike utilisation data and its analysis, several important business decisions will be made, which will be appropriately secured. Data and visualisation reports are only available to a certain group of clients/users.
NFR-3	Reliability	This research offers a trustworthy and effective way to understand how well this bike-sharing programme performed in 2018. Utilizing the IBM Cognos Platform ensures operational report production, upkeep, and accessibility with industry-standard reliability (dashboard).
NFR-4	Performance	The effectiveness of a bike-sharing system in terms of both its spatial and operational efficiency. In order to increase the operational effectiveness of the bike-sharing system, it is critical to assess the state of bike lanes from the viewpoint of public bike riders. The characteristics of bike stations and the distance between bike stations and other amenities are examined by the bike-sharing system dashboard. The evaluation findings can be used to enhance the public bike-sharing service.
NFR-5	Availability	The bicycle-sharing programme is a form of shared transportation in which people can rent bicycles at a reasonable cost for a limited amount of time. CitiBike offers two different kinds of docking systems: docking systems, which allow customers to borrow a bike from one dock and return it to another port within the system; and dockless systems, which are node-free and depend on smart technology. Both forms can use smartphone online mapping to find close-by ports and bikes that are available.
NFR-6	Scalability	Urban inhabitants can immediately get access to bike-sharing programmes, which may make the transportation system more dependable. The programme can be expanded to include locations that are now unreachable by this type of transportation, as well as cities other than New York City, if the necessary data is available and obtained. This research will eventually be able to give a more in-depth picture of how bike-sharing functions in emergency situations as additional data becomes available, particularly in other cities with comparable extensive bike-sharing systems.

3. PROJECT DESIGN

3.1 Data Flow Diagram



User Stories

Use the below template to list all the user stories for the product.

User Type	Functional Requirement (Epic)	User Story Number	User Story / Task	Acceptance criteria	Priority	Release
-----------	-------------------------------	-------------------	-------------------	---------------------	----------	---------

Customer, Analysts, Organizations, Government	Collection of user data	USN-1	Lyft citi bike's official website provides the data to help with analysis, development, visualization etc. Data is collected from these published files.	I can access the data on Lyft citi bike's official website	High	Sprint-1
	Analysing the user data	USN-2	This data is used as input for creating various types of visualizations and analysis is done and a dashboard is created	I can view the analysis of the citi bike	High	Sprint-1
	Dashboard	USN-3	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.	I can register & access the dashboard with login	High	Sprint-2

3.2 Solution & Technical Architecture

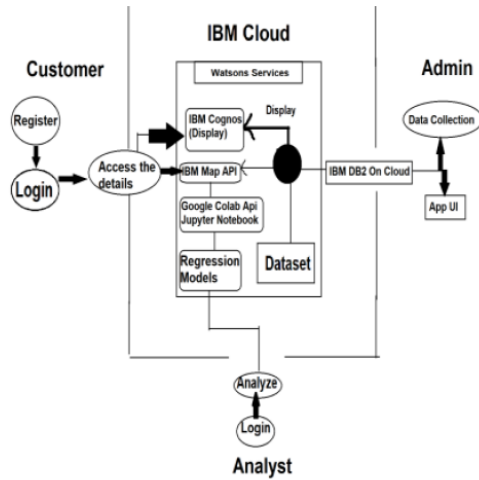
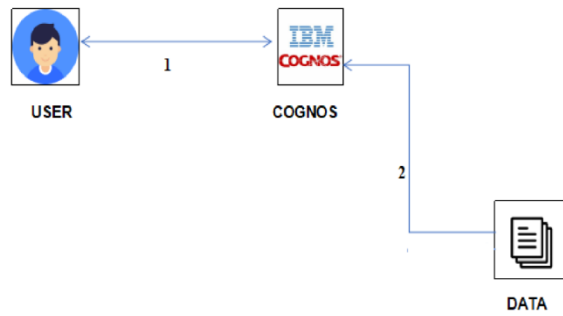
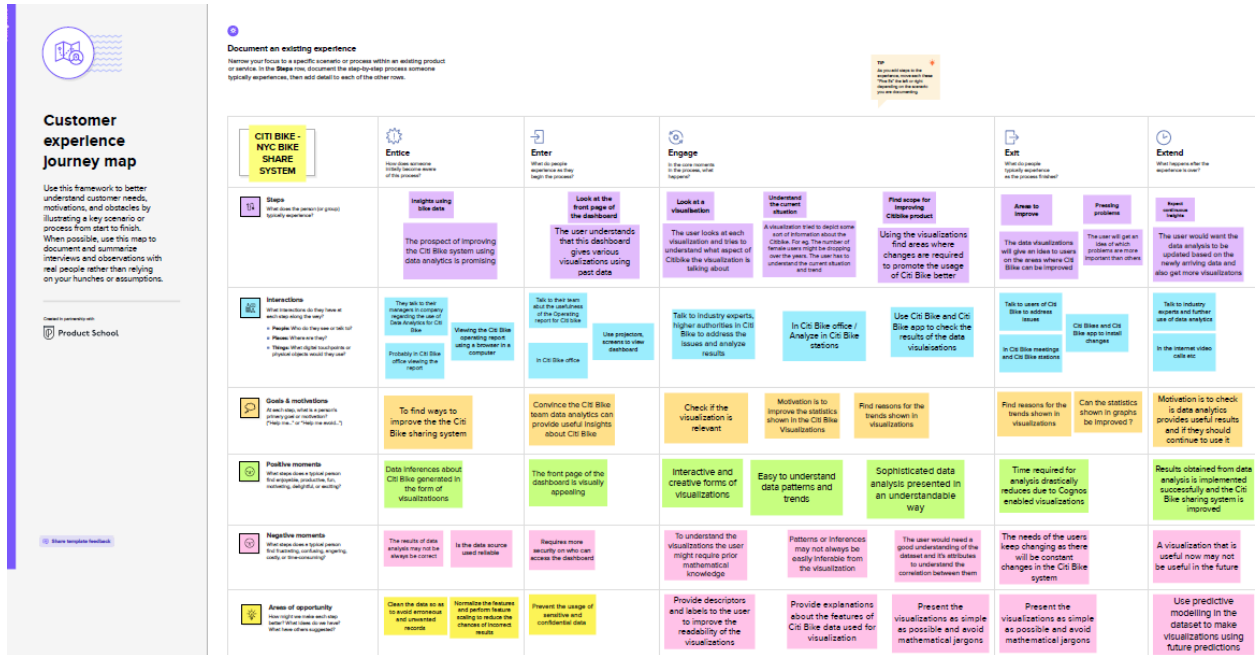


Table-1 : Components & Technologies:

S.No	Component	Description	Technology
	User Interface	1. Display the visualization of the analysed data 2. Display the inferences from the analysed data	HTML, CSS, JavaScript and IBM Cognos

	Application Logic-1	Display details	HTML
	Database	Data Type, Configurations etc.	MySQL
	Cloud Database	Database Service on Cloud	IBM DB2, IBM Cloudant etc.
	File Storage	File storage requirements	IBM Block Storage or Other Storage Service or Local Filesystem
	External API-1	To map the Citi Bike ride in NYC	IBM Map API, etc.
	External API-2	Analysis of the data	Google Colab, Jupyter Notebook
	Machine Learning Model	To plot graphs and predict values	Regression models
	Infrastructure (Server / Cloud)	Application Deployment on Local System / Cloud Local Server Configuration: Local Syst Cloud Server Configuration : IBM Cloud	Local, Cloud Foundry, Kubernetes, etc.

3.3 User Stories



4. PROJECT PLANNING & SCHEDULING

4.1 Sprint Planning & Estimation

Product Backlog, Sprint Schedule, and Estimation (4 Marks)

Use the below template to create product backlog and sprint schedule

Sprint	Functional Requirement (Epic)	User Story Number	User Story/Task	Story Points	Priority	Team Members
Sprint-1	Registration	USN-1	As a user, I can register for the application by entering my email, password, and confirming my password	2	High	T.Ajay Hermas, P.Muthu
Sprint-1		USN-2	As a user, I will receive confirmation email once I have registered for the application	2	High	T.Ajay Hermas, J.Ajay
Sprint-1		USN-3	As a user, I can register for the application through Gmail	2	Medium	T.Ajay Hermas, L.Sumang
Sprint-2	Login	USN-4	As a user, I can log into the application by entering my password	2	High	T.Ajay Hermas, L.Sumang

Sprint	Functional Requirement(Epic)	User Story/Task	Story Points	Priority	Team Members
Sprint-2	Collection of user data	USN-5 I can access and collect the citi bike share system data from Lyft citi bike's official website that has the published files.	2	Medium	T.Ajay Hermas,
Sprint-2		USN-6 I can use the citi bike share system data for analysis purposes	5	High	Ushumgn, T.Ajay Hermas,
Sprint-3	Analysing the user data	USN-7 The data is used as input for creating various types of dashboard to be used for the analysis of the citi bike	8	High	T.Ajay Hermas, P.Muthu J.Ajay Ushumgn
Sprint-3	Dashboard	USN-8 I can register & access the dashboard created for the analysis by logging in	3	Medium	T.Ajay Hermas, P.Muthu
Sprint-3		USN-9 As a user I can view the dashboard that displays the top bike used with respect to trip duration	5	High	Ushumgn
Sprint-4		USN-10 As a user I can view the dashboard that displays the top bike used with respect to customer age group	5	High	T.Ajay Hermas
Sprint-4		USN-11 As a user I can view the dashboard that displays the top bike used with respect to gender	5	High	P.Muthu
Sprint-4		USN-12 As a user I can view the dashboard that displays the total number of trips	5	High	J.Ajay

Sprint	Total Story Points	Duration	Sprint Start Date	Sprint End Date (Planned)	Story Points Completed (as on Planned End Date)	Sprint Release Date (Actual)
Sprint-1	20	6 Days	24 Oct 2022	29 Oct 2022	20	29 Oct 2022
Sprint-2	20	6 Days	31 Oct 2022	05 Nov 2022	20	05 Nov 2022
Sprint-3	20	6 Days	07 Nov 2022	12 Nov 2022	20	12 Nov 2022
Sprint-4	20	6 Days	14 Nov 2022	19 Nov 2022	20	19 Nov 2022

4.2 Sprint Delivery Schedule



5. WORKING WITH THE DATASET & DATA VISUALISATION

5.1 Understanding the dataset

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	tripduration	starttime	stoptime	start	staticstart	staticstart	staticstart	staticend	statio	endstatio	endstatio	endstatio	bikeid	usertype	birth	year	gender		
2	695	01-06-2013 00:00:01	01-06-2013 00:00:01	444 Broadway	40.74235	-73.9892			434 9 Ave & W	40.74317	-74.0037		19678	Subscriber	1983	1			
3	693	01-06-2013 00:00:01	01-06-2013 00:00:01	444 Broadway	40.74235	-73.9892			434 9 Ave & W	40.74317	-74.0037		16649	Subscriber	1984	1			
4	2059	01-06-2013 00:00:01	01-06-2013 00:00:01	406 Hicks St &	40.69513	-73.9996			406 Hicks St &	40.69513	-73.9996		19599	Customer	NULL	0			
5	123	01-06-2013 00:00:01	01-06-2013 00:00:01	475 E 15 St & I	40.73524	-73.9876			262 Washingt	40.69178	-73.9737		16352	Subscriber	1960	1			
6	1521	01-06-2013 00:00:01	01-06-2013 00:00:01	2008 Little West	40.70569	-74.0168			310 State St &	40.68927	-73.9891		15567	Subscriber	1983	1			
7	2028	01-06-2013 00:00:01	01-06-2013 00:00:01	485 W 37 St &	40.75038	-73.9834			406 Hicks St &	40.69513	-73.9996		18445	Customer	NULL	0			
8	2057	01-06-2013 00:00:01	01-06-2013 00:00:01	285 Broadway	40.73455	-73.9907			532 S 5 Pl & S	40.71045	-73.9609		15693	Subscriber	1991	1			
9	369	01-06-2013 00:00:01	01-06-2013 00:00:01	509 9 Ave & W	40.7455	-74.002			521 8 Ave & W	40.75097	-73.9944		16100	Subscriber	1981	1			
10	1829	01-06-2013 00:00:01	01-06-2013 00:00:01	265 Stanton St	40.72229	-73.9915			436 Hancock S	40.68217	-73.954		15234	Subscriber	1984	1			
11	829	01-06-2013 00:00:01	01-06-2013 00:00:01	404 9 Ave & W	40.74058	-74.0055			303 Mercer St	40.72363	-73.9995		16400	Subscriber	1987	1			
12	1316	01-06-2013 00:00:01	01-06-2013 00:00:01	423 W 54 St &	40.70585	-73.9869			314 Cadman P	40.69383	-73.9905		19781	Subscriber	1960	1			
13	1456	01-06-2013 00:00:01	01-06-2013 00:00:01	502 Henry St &	40.71422	-73.9813			532 S 5 Pl & S	40.71045	-73.9609		18886	Customer	NULL	0			
14	386	01-06-2013 00:00:01	01-06-2013 00:00:01	241 DelKalb Av	40.68981	-73.9749			365 Fulton St	40.68223	-73.9615		19039	Subscriber	1981	1			
15	924	01-06-2013 00:00:01	01-06-2013 00:00:01	486 Broadway	40.7462	-73.9886			521 8 Ave & W	40.75097	-73.9944		16608	Customer	NULL	0			
16	1233	01-06-2013 00:00:01	01-06-2013 00:00:01	527 E 33 St &	40.74402	-73.9761			296 Division St	40.71413	-73.997		14761	Subscriber	1987	1			
17	512	01-06-2013 00:00:01	01-06-2013 00:00:01	309 Murray St	40.71498	-74.013			300 Shevcheni	40.72815	-73.9902		19080	Subscriber	1979	2			
18	505	01-06-2013 00:00:01	01-06-2013 00:00:01	309 Murray St	40.71498	-74.013			347 Greenwid	40.72885	-74.0086		16798	Subscriber	1984	1			
19	833	01-06-2013 00:00:01	01-06-2013 00:00:01	503 E 20 St & F	40.73827	-73.9875			503 E 20 St & F	40.73827	-73.9875		19072	Customer	NULL	0			
20	1818	01-06-2013 00:00:01	01-06-2013 00:00:01	257 Lispenard	40.71939	-74.0025			500 Broadway	40.76229	-73.9834		20349	Customer	NULL	0			

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

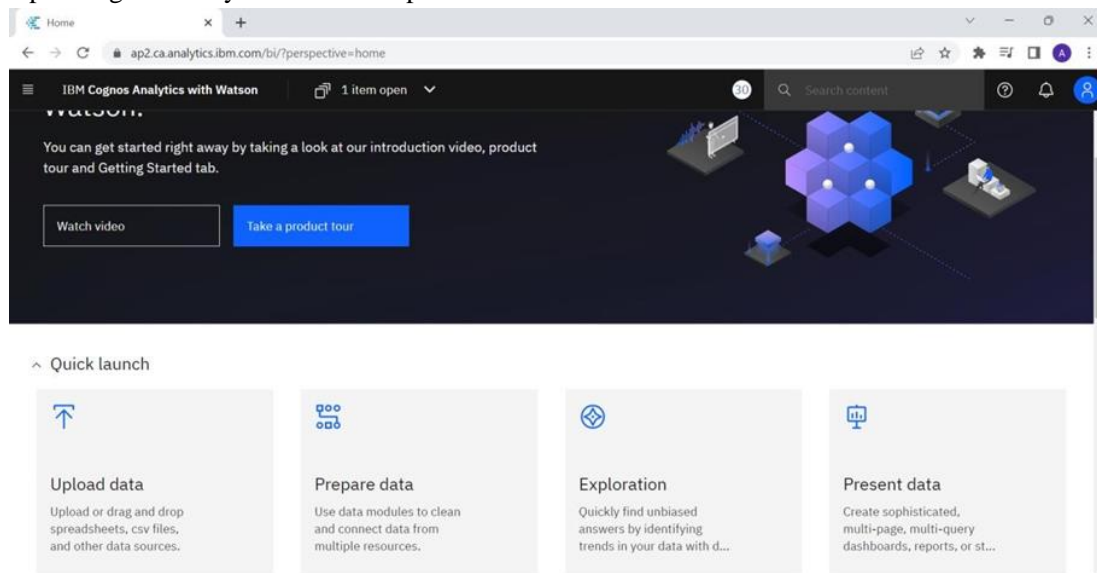
13.User Type (Customer = 24-hour pass or 3-day pass user; Subscriber = Annual Member):
Customers are usually tourists, subscribers are usually NYC residents

14.Year of Birth: Self-entered, not validated by an ID

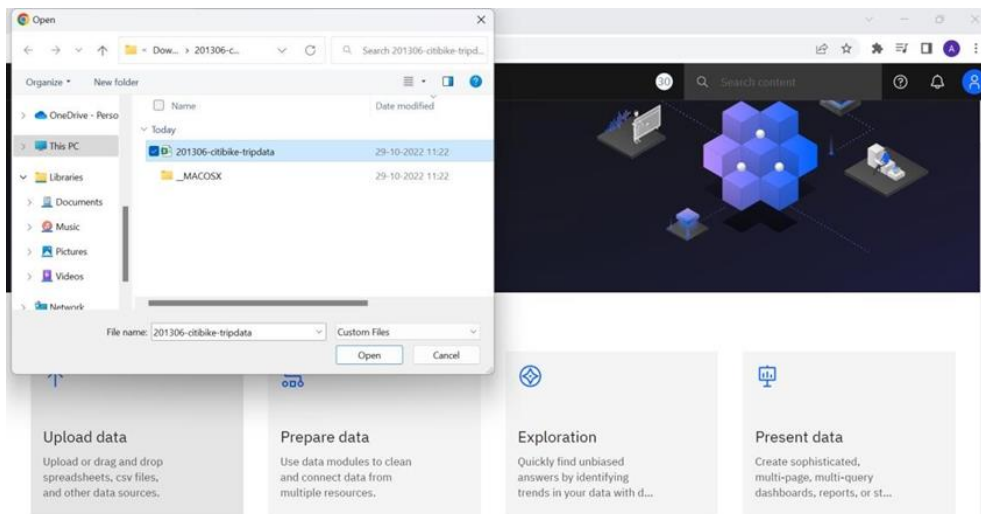
Gender (Zero=unknown; 1=male; 2=female): Usually unknown for customers since they often sign up at a kiosk

7.2 Loading the dataset

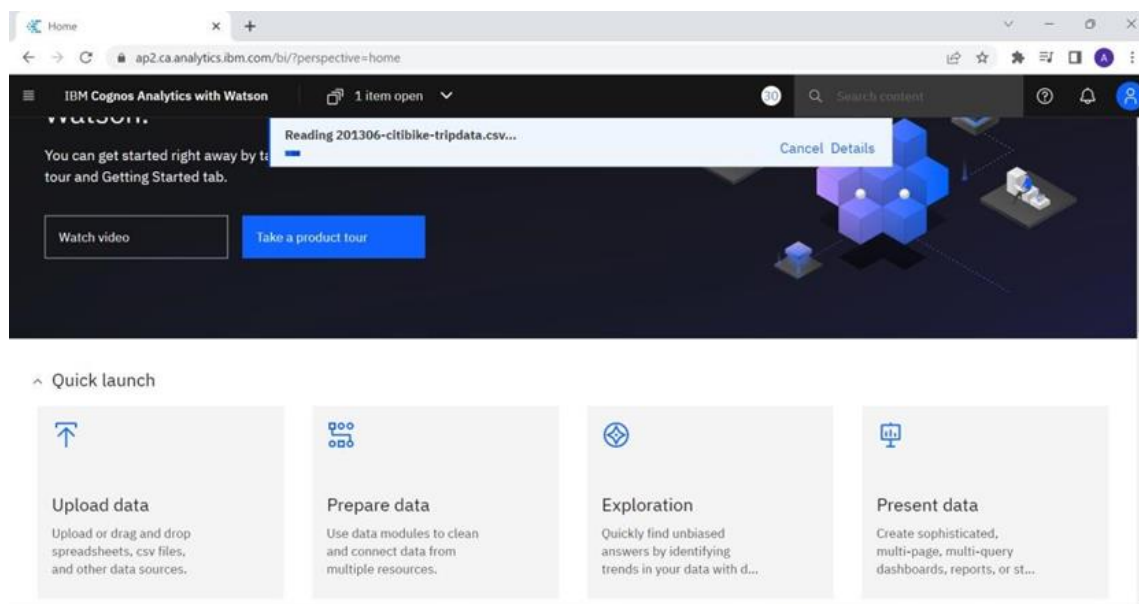
Open Cognos Analytics and click upload data



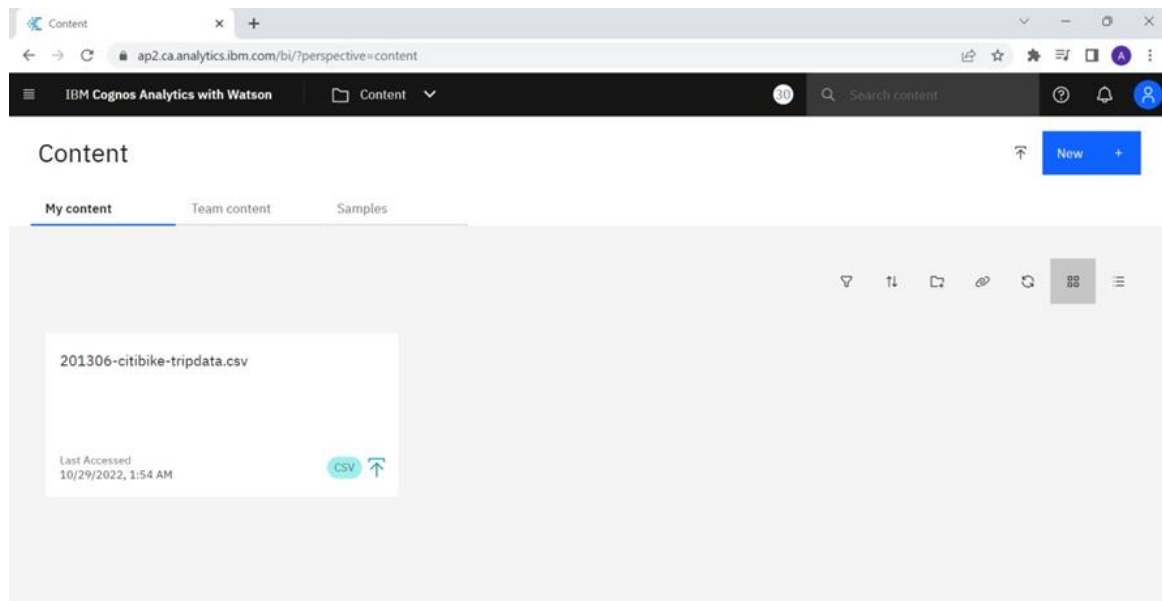
Select the dataset to be uploaded



The excel file is getting uploaded in Cognos Analytics

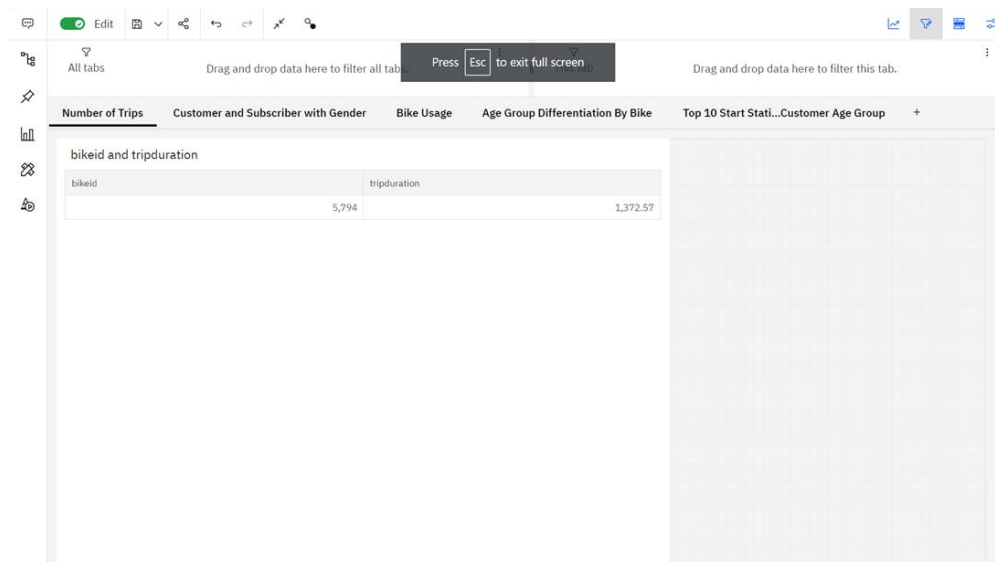


The dataset can be accessed in My Content in Cognos Analytics

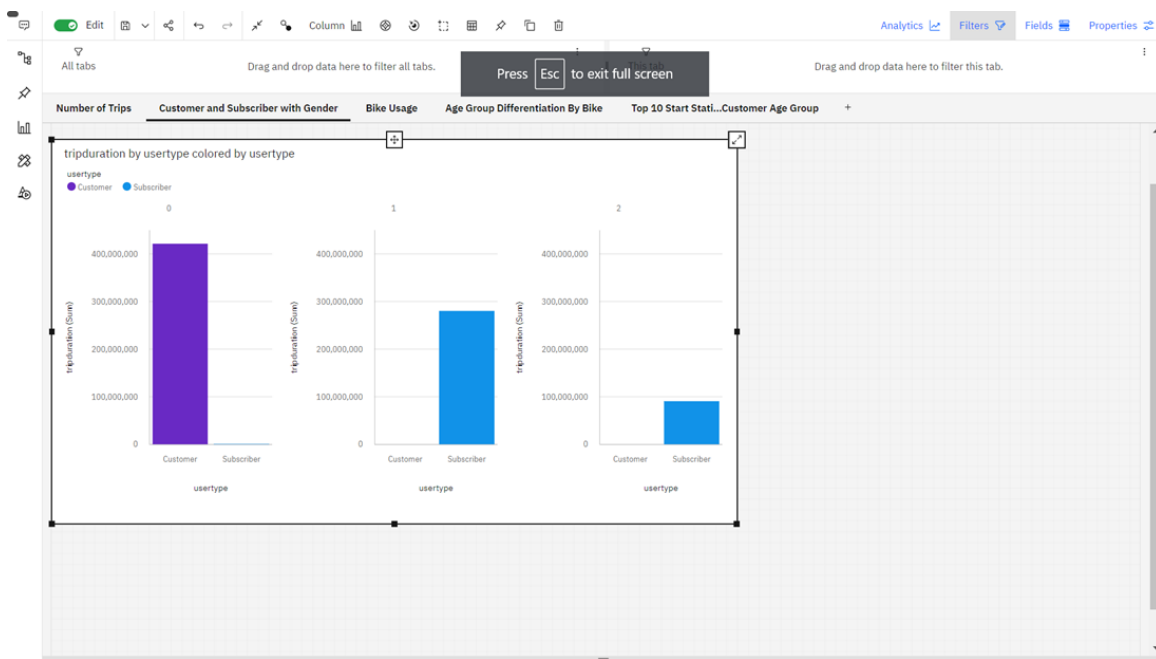


7.3 Visualization charts

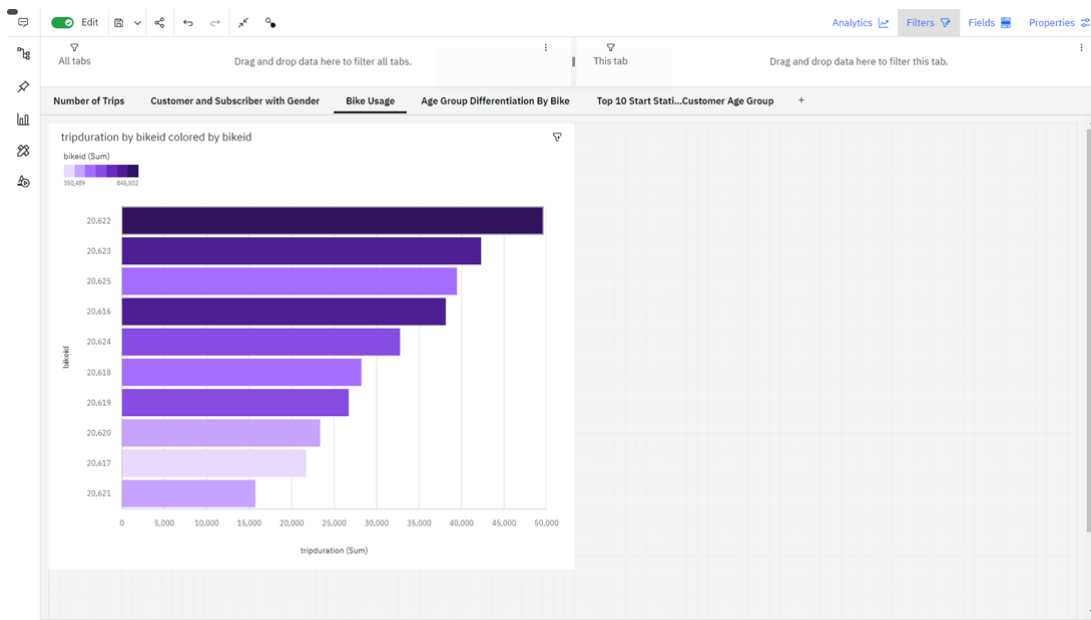
Number of Trips:



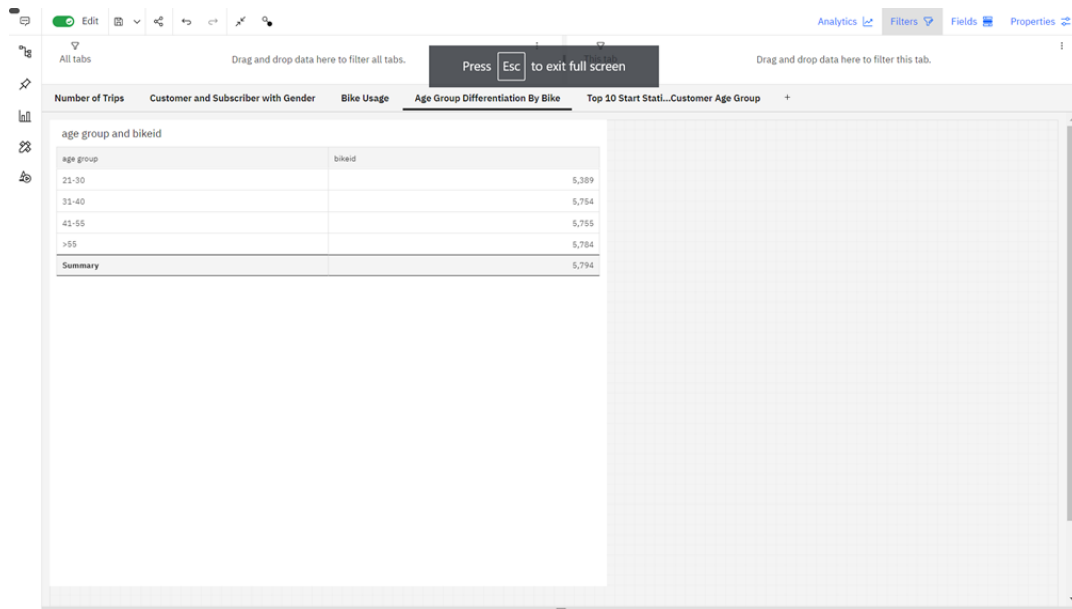
Customer and Subscriber with Gender:



Bike Usage:



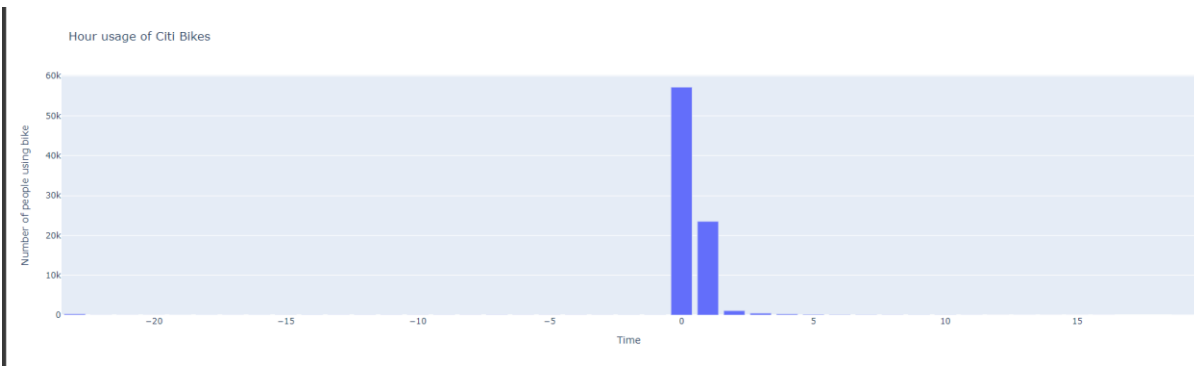
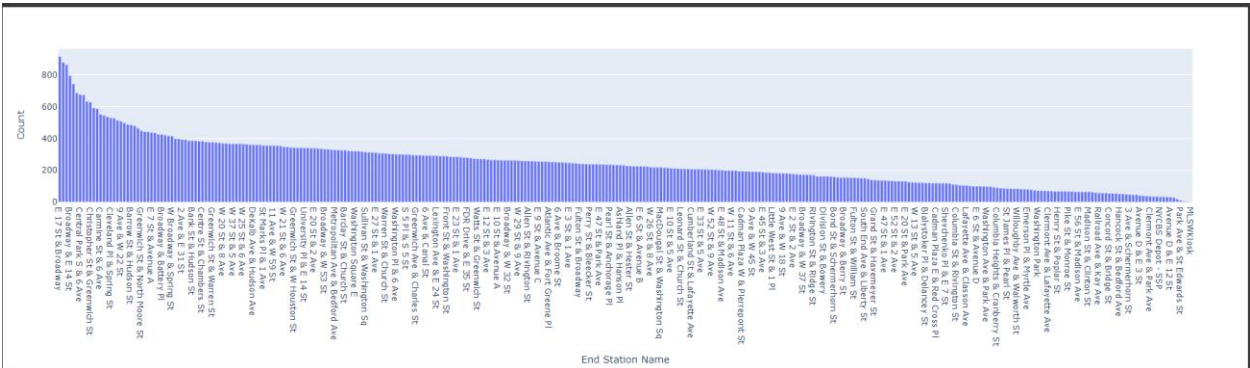
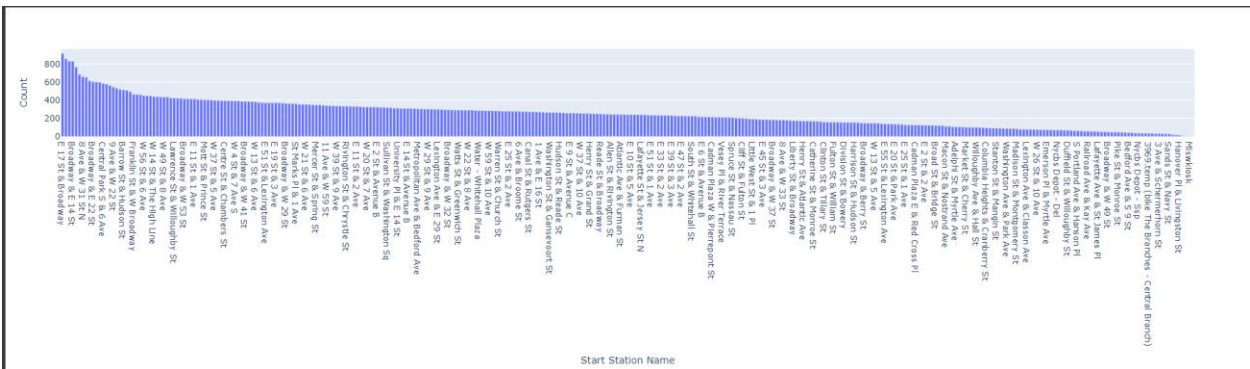
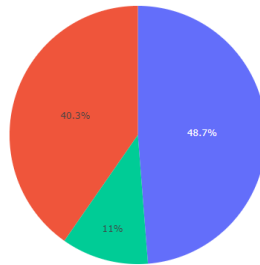
Age group differentiation by bike:



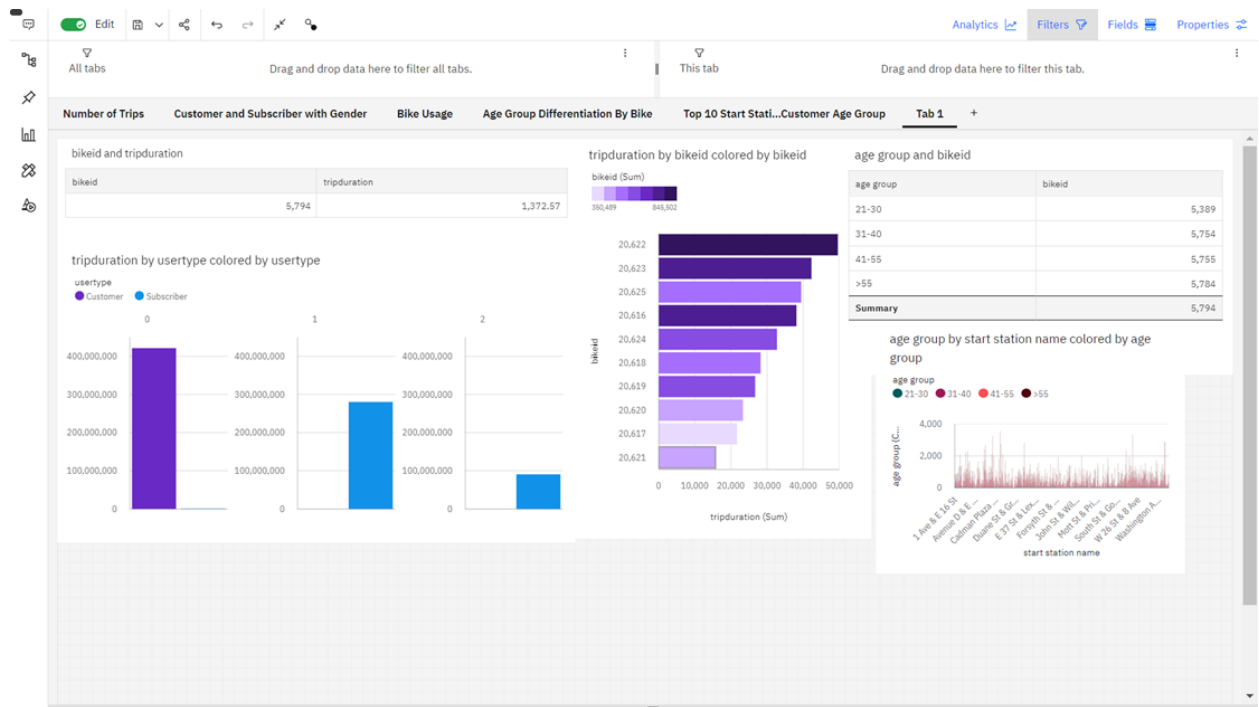
Top 10 Start Station Names with Respect to Customer Age Group:



Gender Variation



6. CREATING THE DASHBOARD



7. 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.

8. 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.

9. 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.

10. 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"})

```

11. GITHUB LINK

<https://github.com/IBM-EPBL/IBM-Project-46101-1660737926>