

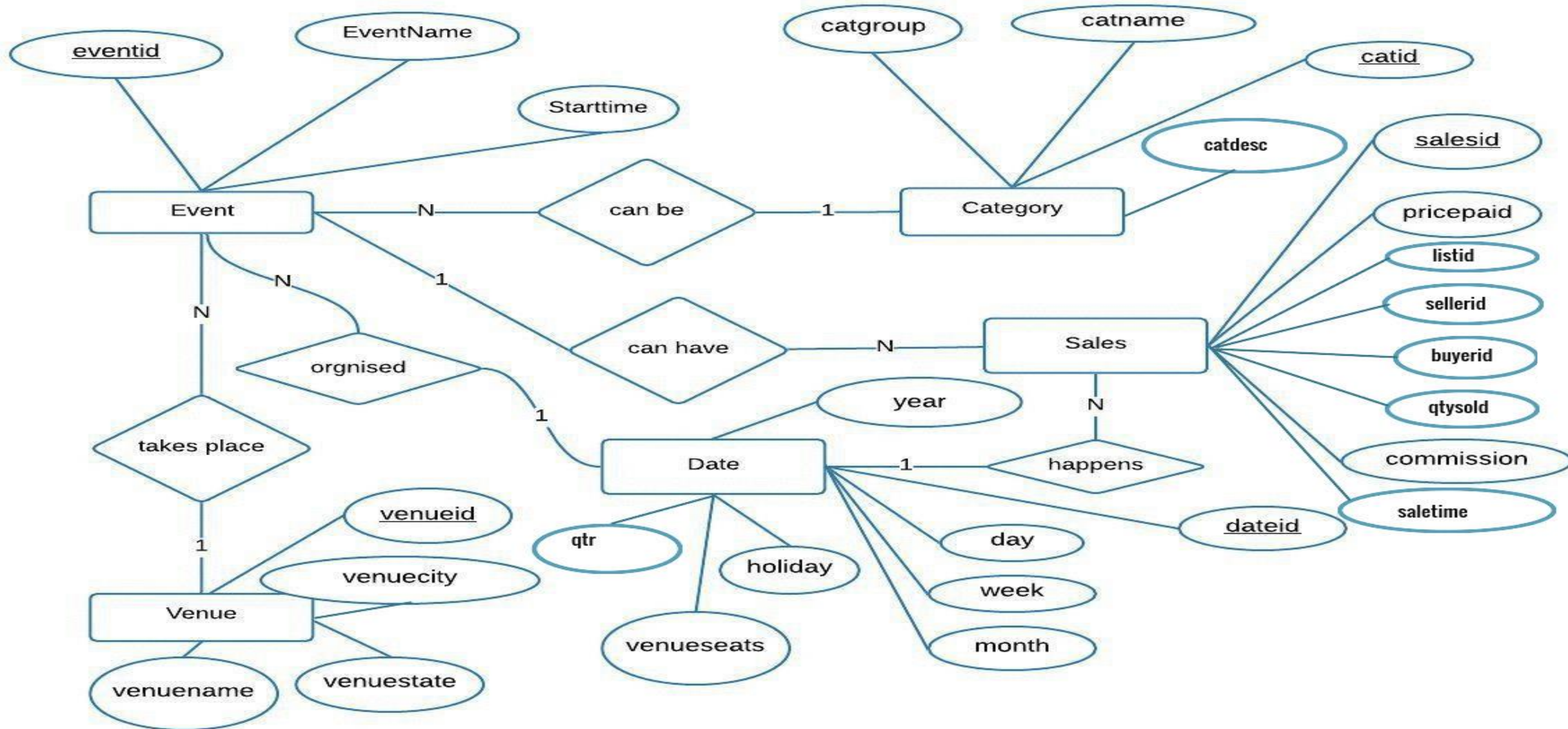


Database System S21

Assignment 10

- Suresh Kumar Choudhary 5504392
- Rohan Chitte 5509596
- Frenny Macwan 5504865

Task 1: ER diagram



Task 1: Relational Schema

- Event (eventid, venueid, catid, dateid, eventname, starttime)
- category(catid, catgroup, catname, catdesc)
- Sales (salesid, eventid, dateid, listid, sellerid, buyerid, qty sold, price paid, commission, sale time)
- Venue
(venueid, venue name, venue city, venue state, venue seats)
- Date (dateid, cal date, day, week, month, qtr, year, holiday)

Objective

- **Analysis of bike sales of year 2008 in the USA**
 - Monthly sales
 - USA state wise sales
 - Events wise sales
 - Category wise sales (ie. Musicals, opera, plays, pop)

Task 2: Cleaning

The dataset contains multiple csv files. So, Pandas library in python was used to read and store the csv files as dataframes.

```
In [1]: 1 import numpy as np
        2 import pandas as pd
```

```
In [2]: 1 #Reading csv files into dataframe
        2 category = pd.read_csv(r'/users/rohanchitte/downloads/data-2/category.csv')
        3 date_tab = pd.read_csv(r'/users/rohanchitte/downloads/data-2/date_table.csv')
        4 event = pd.read_csv(r'/users/rohanchitte/downloads/data-2/event.csv')
        5 sales = pd.read_csv(r'/users/rohanchitte/downloads/data-2/sales.csv')
        6 venue = pd.read_csv(r'/users/rohanchitte/downloads/data-2/venue.csv')
```

```
In [3]: 1 print(venue.isnull().sum())
        2 print(category.isnull().sum())
        3 print(date_tab.isnull().sum())
        4 print(event.isnull().sum())
        5 print(sales.isnull().sum())
```

```
venueid      0
venueid      0
venueid      0
venueid      0
venueid      0
venueid      15
dtype: int64
```

Task 2: Cleaning (Cont.)

Pandas `isnull()` function was used to find null values and only the column `venueseats` of `venue` dataframe had null values.

```
In [ ]: 1 # 130/200 values are 0, so we will replace null values with 0.
        2 #Moreover, the objective of visualization does not deal with venue seats.
        3 #So, data imputation becomes irrelevant.
        4 venueseats = venue['venue seats']
        5 count = 0
        6 for i in venueseats:
        7     if i == 0:
        8         count = count + 1
        9 print(count) #130
       10
```

```
In [ ]: 1 #Data Imputation
        2 venue['venue seats'] = venue['venue seats'].fillna(0)
```

Task 2: Import Data

Postgresql Database was used to create a database. Python package sqlalchemy was used to first establish a connection with the database. Then, Pandas function to_sql() was used to write and store records of the dataframe to the database.

```
In [17]: 1 #Importing Data
```

```
In [26]: 1 from sqlalchemy import create_engine  
2 engine = create_engine('postgresql://:@localhost:5432/databaseproject')
```

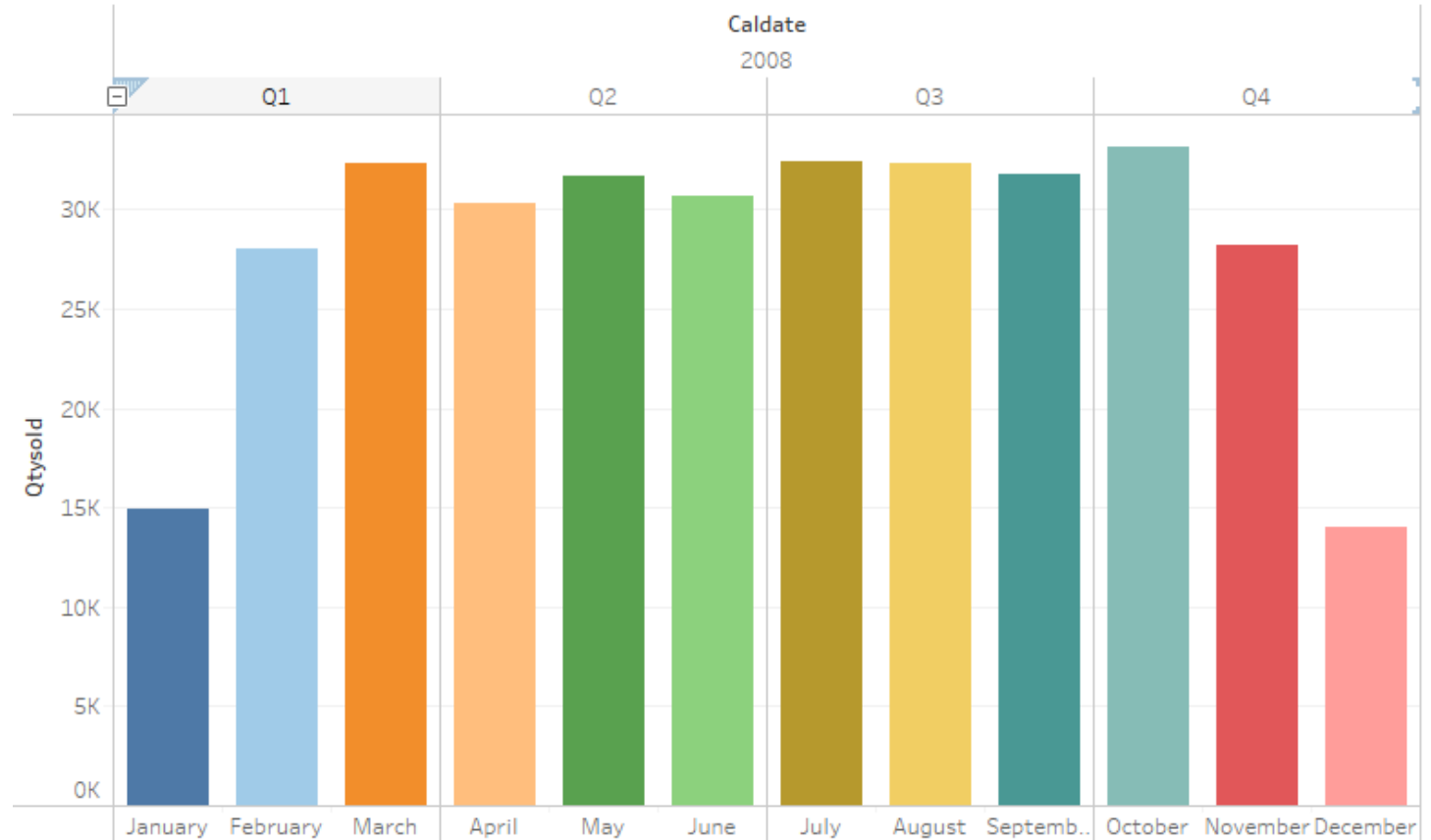
```
In [27]: 1 sales.to_sql('sales', engine)  
2 venue.to_sql('venue', engine)  
3 event.to_sql('event', engine)  
4 date_tab.to_sql('date_table', engine)  
5 category.to_sql('category', engine)
```

Task 3:

Visualization of the data set

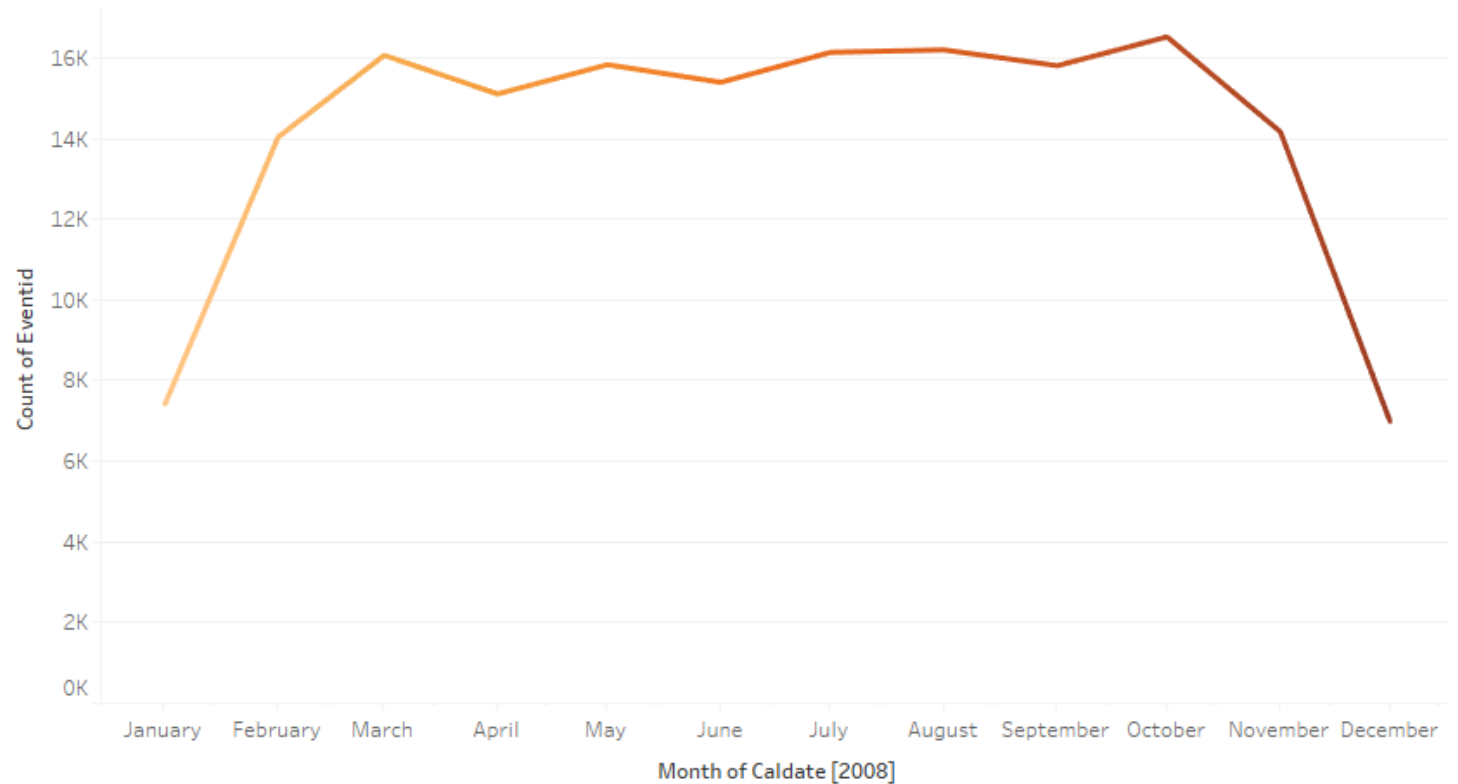
- Analysis of bike sales for year 2008.

What is particularly encouraging is to see how well e-bikes sales have held up in the USA in the challenging market conditions for year 2008. It can be seen from 'Monthly sales' graph, the quantity of bikes is sold is quite steady except for two months – January and December.



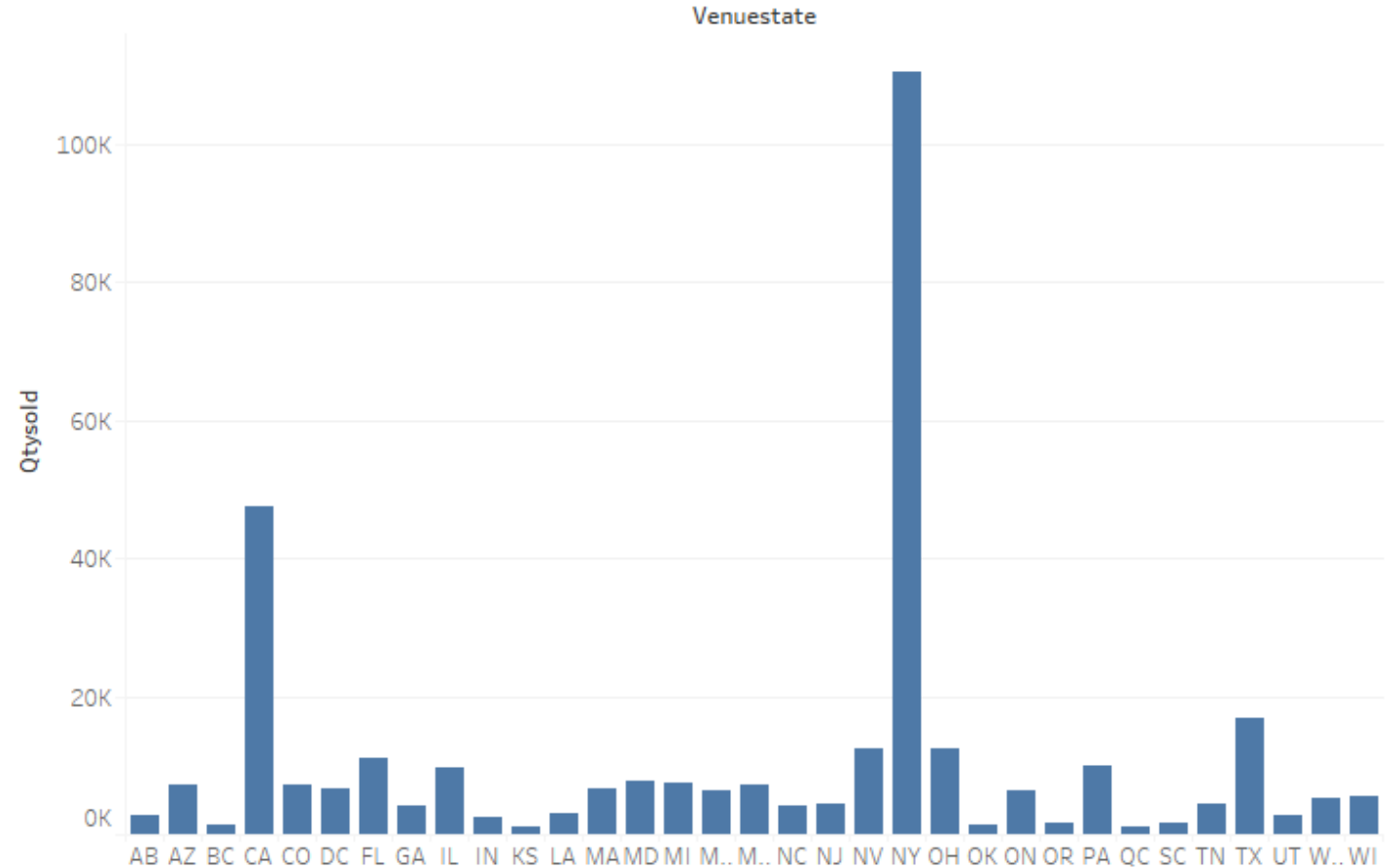
Task 3: Visualization... Cont...

From above graph, number of events happened in January and December months are less compared to other months. That also reflects into quarterly sells. So, we can say that sales of bikes have been increased when more events are happening throughout the year.



Task 3: Visualization... Cont...

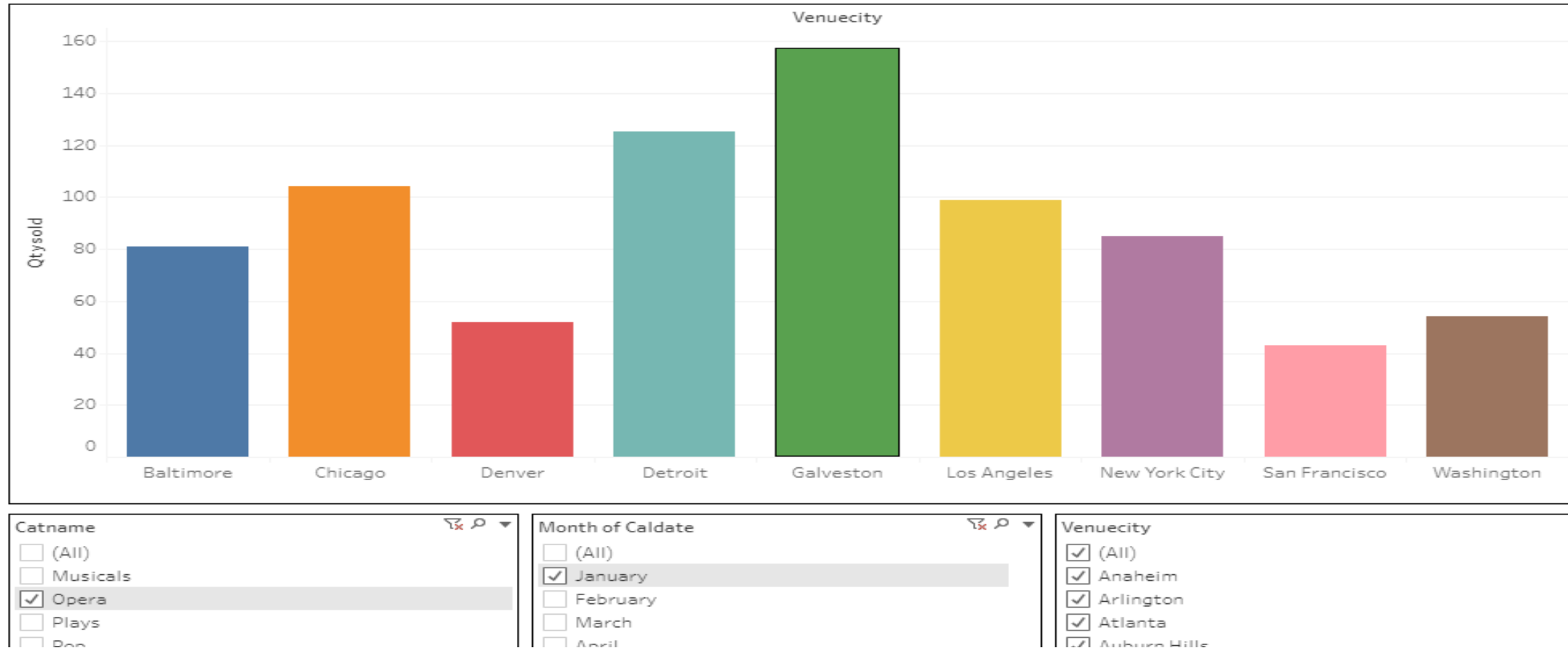
This gives a general overview of the sale of e-bikes in different states of USA. However, it also confirms some known challenges with the different states of the country. We already know NY and CA are main states of the USA, where a greater number of events are taking place. It can be inferred that in these two states, more numbers of bikes are purchased compared to other states.



Task 3: Visualization...

Cont...

Analysis of sales in different US cities



- Source code:



Assignment10.html



C:\Users\Frenny
ran\Downloads\Da

- Github link : <https://github.com/FrennyMacwan/Database-Project>

*Thank
you!*