

**MISM6205 – DATA WRANGLING FOR BUSINESS**  
**PROJECT REPORT**

<b>Group Number</b>	6
<b>Group Members</b>	Ayushi Anantwar (002751798) Pooja Pendharkar (002653116) Ramya Kumaresan (002990437)
<b>Project Goal</b>	To recommend must-visit destinations in proximity to major attractions in New York City based on the type of tourists.

**TABLE OF CONTENTS**

<b><u>CONTENT</u></b>	<b><u>PAGE NUMBER</u></b>
TOPIC AND BUSINESS QUESTION	<b>2</b>
DATASET AND DATASOURCE DETAILS	<b>3</b>
INFORMATION QUALITY	<b>4</b>
METHODS AND TOOLS	<b>5</b>
DATA WRANGLING PROCESS	<b>6</b>
ANALYSIS AND RESULTS	<b>7</b>
EXTERNAL ARTICLES	<b>11</b>
ADDITIONAL DETAILS AND ANALYSIS	<b>11</b>

## TOPIC AND BUSINESS QUESTIONS



This project aims to recommend must-visit destinations in proximity to some major attractions in New York City based on the type of tourists. The business-related question that we are planning to explore through this project is - “How to plan a short yet fulfilling trip for the incoming tourists”. Our project aims to maximize the number of locations tourists can visit within a short period. It is assumed that tourists will visit important attractions in the city such as the Statue of Liberty, Central Park, Times Square, and the Empire State Building. So, the results of this project will aid tourists in identifying some interesting places such as eateries, museums, art galleries, beaches, etc. around the four major attractions in New York City.

New York City is taken as a use case in this project as it is popularly known as the world’s busiest tourist destination. Over 60 million tourists visit New York City every year, of this number 12 million come from outside the United States. Most of these tourists depend on travel blogs or different websites to gather information on the best places to visit.

The trip planners available in the market are based on travel dates or the number of days for the trip and showcase a list of top places to visit as a recommendation. However, the restrictions of travel time between these places are not considered. The algorithms do not recommend a list that would allow the tourist to visit multiple locations in a short amount of time. For example, Wanderlog recommends a generic list of places to visit. Roadtrippers shows a list of tourist attractions based on route and considers the travel time, however, it is focused specifically on road trips i.e., intercity travel instead of intracity. Most tourists plan their trips during weekends, office leaves, or holidays. They wish to plan a trip that helps them make the most of the available time i.e., being able to cover as many attractions as possible. This project aims to plan such a trip for tourists by recommending must-visit places around the top attractions that one can cover according to their interests.

What makes our project unique is that it not only suggests the must-visit places based on the location proximity but also considers the type of tourists visiting New York City (e.g.: Individuals, Friends, Family, and Couples). This feature of our project helps tourists to choose what is most likely to suit and accommodate their travel preferences from a range of options.

Through our project, we wish to take a step in this direction to help tourists plan a trip with an optimized list of attractions they can visit. The recommendations are shortlisted by traversing multiple datasets to showcase the places which are closest to a specific tourist attraction.

### **DATASET AND DATA SOURCE DETAILS**

The data collection process required us to collect the below-listed datasets that contain most of the required information about the tourist attractions and eateries in New York City. The “**ZIPCODE**” would be used to recommend the best attractions in and around the specific area of choice entered by the user. We have added a separate column called “**CATEGORY**” in each of the datasets to specify the category for each tourist attraction.

<b>Dataset Name</b>	<b>Dataset URL</b>	<b>Variables used for our analysis</b>	<b>Dataset Provided By</b>
<b>Museums</b>	<a href="https://data.cityofnewyork.us/Recreation/New-York-City-Museums/ekax-ky3z">https://data.cityofnewyork.us/Recreation/New-York-City-Museums/ekax-ky3z</a>	<b>Name</b> - name of museum, <b>Tel</b> - phone no, <b>URL</b> - website, <b>Address1</b> - Street name, <b>Address2</b> - building name or floor, <b>Zip</b> – zipcode	Department of Information Technology & Telecommunications (DoITT)
<b>Art Galleries</b>	<a href="https://data.world/city-of-ny/tgyc-r5jh">https://data.world/city-of-ny/tgyc-r5jh</a>	<b>Name</b> - Name of the art gallery, <b>Tel</b> - Contact Number of the art gallery, <b>URL</b> - Link to the art gallery's website, <b>Address1</b> of the art gallery, <b>Address2</b> - Additional details about address of the art gallery, <b>Zip</b> - Zipcode of the area where the art gallery is located	Data World
<b>Beaches</b>	<a href="https://data.world/city-of-ny/zyf6-z3xt">https://data.world/city-of-ny/zyf6-z3xt</a>	<b>Name</b> -Name of the property, <b>location</b> - gives the location details of the beach, <b>Phone</b> -Contact number, <b>Lat and long</b> - shows the latitudes and longitudes of beaches, <b>Description</b> – Contains a few words about the beaches	Data World
<b>Botanical Gardens</b>	<a href="https://data.world/city-of-ny/hrii-hezi/workspace/file?filename=botanical-gardens-1.csv">https://data.world/city-of-ny/hrii-hezi/workspace/file?filename=botanical-gardens-1.csv</a>	<b>Name</b> - different names of the botanical garden, <b>Tel</b> - phone number of different botanical gardens, <b>URL</b> - website of the different botanical gardens, <b>Address1</b> – it shows the address of different botanical gardens in NY,	Data World

		<b>Zipcode</b> -different area zip codes of different botanical gardens	
<b>Water Trails</b>	<a href="https://data.world/city-of-nyc/hxay-3qcw/workspace/file?filename=new-york-city-water-trail-kayak-and-canoe-launch-sites-1.json">https://data.world/city-of-nyc/hxay-3qcw/workspace/file?filename=new-york-city-water-trail-kayak-and-canoe-launch-sites-1.json</a>	<b>Name</b> - names of the water trails and sport, <b>Location</b> - locations of water trails and sport, <b>Description</b> - more information about water trails and sport, <b>Lat and long</b> - shows the latitudes and longitudes of water trails and sport	Data World
<b>Hiking Trails</b>	<a href="https://data.world/city-of-nyc/i8f4-bu5r/workspace/file?filename=directory-of-hiking-trails-1.json">https://data.world/city-of-nyc/i8f4-bu5r/workspace/file?filename=directory-of-hiking-trails-1.json</a>	<b>Name</b> - Name of the hiking trail, <b>location</b> - location of the hiking trail, <b>Other_Details</b> -Additional notes	Data World
<b>Restaurants</b>	<a href="https://data.cityofnewyork.us/Health/Restaurants-rolled-up-/59dk-tdhz">https://data.cityofnewyork.us/Health/Restaurants-rolled-up-/59dk-tdhz</a>	<b>DBA</b> -Name of restaurants, <b>BORO</b> -boroughs of New York City, <b>Street</b> -Street where the restaurant is located, <b>Zipcode</b> - Zipcode	Department of Health and Mental Hygiene (DOHMH)
<b>Zip code Proximity Data - Manually gathered from the website</b>	<a href="https://ny.postcodebase.com/zipcode radius#myarticle">https://ny.postcodebase.com/zipcode radius#myarticle</a>	<b>ZIP Code_Statue</b> - Zip code of location close to Statue of Liberty, <b>Distance_Statue</b> - distance of location zip code from the statue of liberty, <b>ZIP Code_Empire</b> - Zip code of location close to Empire State Building, <b>Distance_Empire</b> - distance of location zip code from Empire state building, <b>ZIP Code_Times</b> - Zip code of location close to Times Square, <b>Distance_Times</b> - distance of location zip code from Times Square, <b>ZIP Code_Central</b> - Zip code of location close to Central Park, <b>Distance_Central</b> -distance of location zip code from Central Park	Manually collected from website - <a href="https://ny.postcodebase.com/zipcode radius#myarticle">https://ny.postcodebase.com/zipcode radius#myarticle</a>

### **INFORMATION QUALITY**

- The addresses for a few datasets were not separated as Address and Zip code, this made it difficult for us to extract the data and add it to the appropriate columns. This issue was resolved during the data-wrangling process.
- Since we did not have any readily available dataset to determine the zip code proximity for the four major tourist attractions (Statue of Liberty, Empire State Building, Central

Park, New York Times Square) in New York, the zip code proximity dataset for each of the four major tourist attractions had to be manually created to further enrich the data used for our analysis.

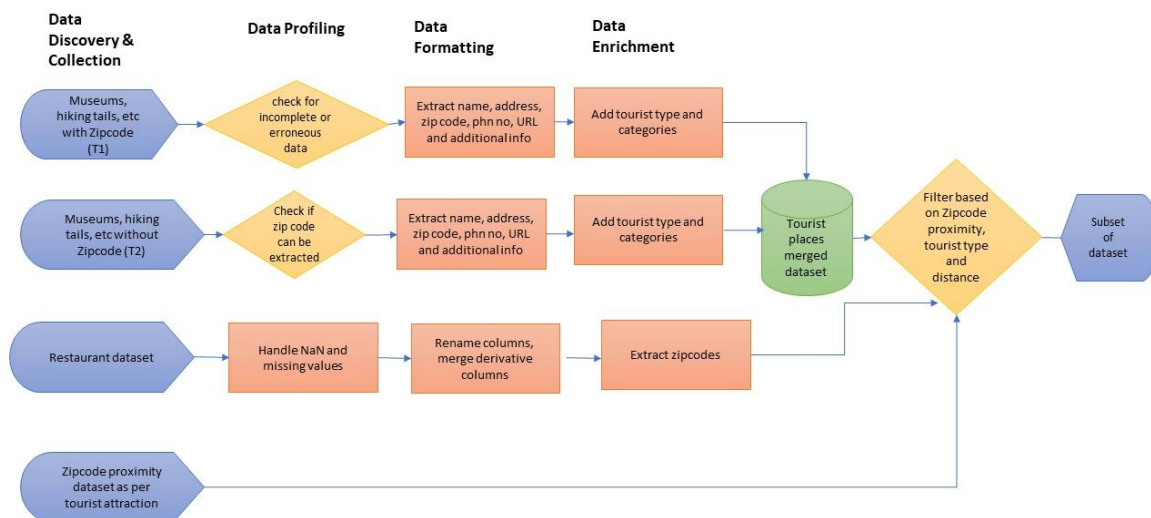
- The process or methodology of data collection is not provided for the datasets from the website “Data World” which makes the reputability of the dataset questionable. However, we went ahead with using these datasets by manually cross-verifying some of the information from the datasets that would be used for our analysis.

## METHODS AND TOOLS

Listed below are the common methods used for wrangling the data post data collection -

- **Data Discovery and Collection** – The datasets required to perform analysis and solve the business problem had to be collected and cleaned to suit the business requirements.
- **Data Profiling** – The datasets had to be analyzed for understanding the purpose and validity of different values and attributes.
- **Data Formatting** – The datasets were formatted to suit the requirements of the analysis.
- **Data Enrichment** – New columns were added to the existing datasets and a new dataset was manually created to help enrich the information provided in the originally collected datasets.

The below diagram illustrates the data-wrangling concepts used for our project -



Tools -

The tools that were used for our project include:

- Microsoft Excel

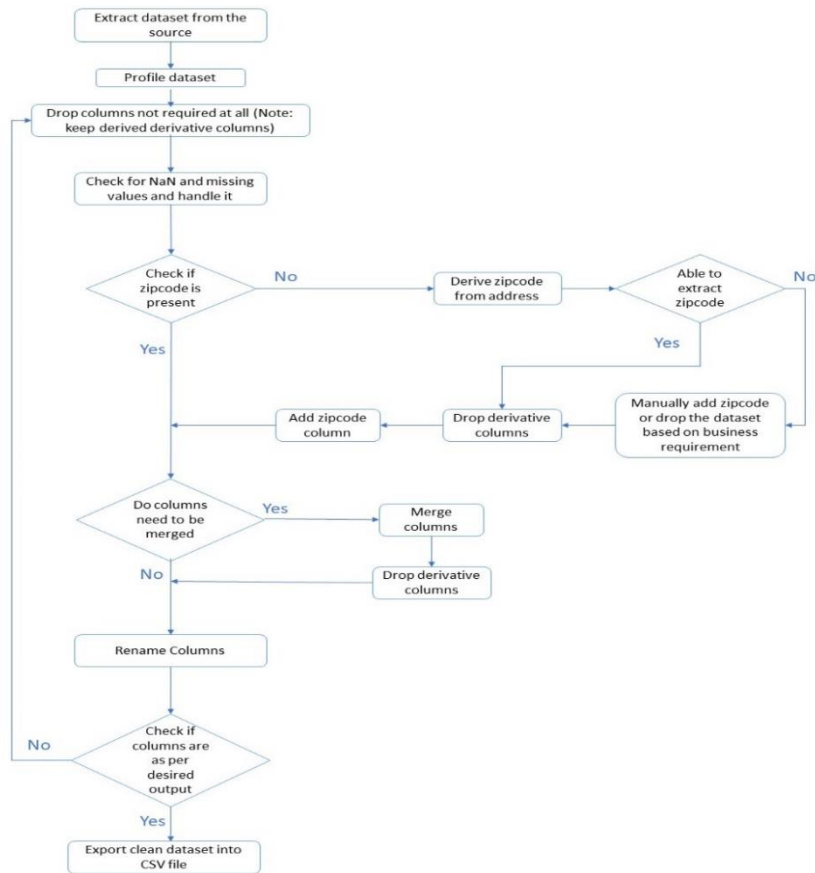
- Python - Pandas, Geopy, PySimpleGUI
- Tableau

#### Challenges -

- The process of fetching the data on zip code proximity for the four major tourist attractions in New York was a challenge during the data collection phase.
- Generating the zip codes from the “**ADDRESS**” field for those tourist attraction datasets that did not have a value for zip code was a challenge. To ease this process, we used the Geopy library in python. However, we were unable to generate the zip code values for certain locations as Geopy could not detect certain tourist locations based on the address parameter that was passed. Instead of returning the correct zip code values, it returned NaN values. To avoid this issue, we had to manually enter the zip codes for a few datasets.
- This project required us to use the PySimpleGUI library for the first time, it took a while for us to learn about the library and implement its functionalities in our project. Making use of this library allowed us to establish a better user-system interaction for the project to provide personalized recommendations.

### **DATA WRANGLING PROCESS**

The diagram below depicts the data-wrangling process used for our project -



Some of the validation rules and checks implemented for our project are as follows -

- Checked the datatype of each column in the individual datasets to ensure uniformity in the datatypes of the column values before merging the datasets to perform the required analysis.
- Formatted the columns (renamed, arranged them in the correct sequence, dropped the derivative columns after merging them to get the required data value) in each of the datasets to ensure overall uniformity in columns before merging the datasets.
- Ensured that each of the files was in the right format before merging them to obtain the final dataset for tourist attractions.
- Checked if the “**ZIPCODE**” column existed for each of the individual datasets. If no, manually added the zip code values.
- Handled the missing and the null values in each of the individual datasets before merging.

## ANALYSIS AND RESULTS

Final Output -

Post fetching the required datasets and wrangling them to obtain the required information. We performed an analysis to determine the secondary tourist attractions and restaurants situated close to the four most popular and primary tourist spots in New York which include the Statue of Liberty, Central Park, New York Times Square, and the Empire State Building based on two primary filters: Distance in miles and tourist type.

Attached below are the output screenshots –

Step 1: The user chooses one of the primary tourist attractions (Major Tourist Attraction) from the drop-down list



Step 2: The user chooses the distance (in miles) between 1 to 5 from the drop-down list to determine the secondary tourist attractions situated within the specified distance (in miles) from the primary tourist attraction/Major Tourist Attraction





Step 3: The user chooses the Tourist Type from the drop-down list to get the personalized recommendation



Welcome!

Tell us about your preferences and we will recommend the best places to visit to make your trip unforgettable!

Major Attraction:  
Empire State Building

Distance (in miles):  
2

Tourist Type:  
Individual  
Couple  
Family  
Friends

CANCEL

Step 4: The final output is printed in the form of two tables:

- The secondary tourist attractions within the specified distance range for the chosen tourist type
- The restaurants present within the specified distance range
- Each of the above-mentioned tables contains the following details: "**CATEGORY**" of the tourist attraction/ the "**CATEGORY**"- Restaurant, "**NAME**" of the tourist attraction, Exact "**ADDRESS**" of the tourist attraction/Restaurant, Exact "**ZIPCODE**" of the tourist attraction/Restaurant, "**PHONE NO.**" (If available), "**URL**" (if available) and some "**ADDITIONAL DETAILS**" about that tourist spot/restaurant.

Recommended Places to visit								
Here are the places you can visit:								
CATEGORY	NAME	ADDRESS	ZIPCODE	PHONE NO.	URL	ADDITIONAL D	Distance	Empi
Art Gallery	Brenda Taylor	511 W. 25th St	10001	(212) 463-7161	http://www.bre	Not Available	0.24 Miles	
Art Gallery	Caelum Galler	526 W 26th St	10001	(212) 924-4161	http://www.cae	Not Available	0.24 Miles	
Art Gallery	Century Artist	530 W 25th St	10001	(212) 367-7071	http://www.nev	Not Available	0.24 Miles	
Art Gallery	Ceres Gallery	547 W 27th St	10001	(212) 947-6101	http://www.cer	Not Available	0.24 Miles	
Art Gallery	Chappell Galle	526 W 26th St	10001	(212) 414-2671	http://www.chi	Not Available	0.24 Miles	
Art Gallery	Cheim & Reac	547 W 25th St	10001	(212) 242-7721	http://www.chi	Not Available	0.24 Miles	
Art Gallery	Clampart	531 W 25th St	10001	(646) 230-0021	http://www.nile	Not Available	0.24 Miles	
Art Gallery	Clementine G	623 W 27th St	10001	(212) 243-5931	http://www.cle	Not Available	0.24 Miles	
Art Gallery	Coploff Gallery	526 W. 26th St	10001	(212) 674-1021	http://www.arti	Not Available	0.24 Miles	
Art Gallery	Bowery Galler	530 W 25th St	10001	(646) 230-6651	http://www.boi	Not Available	0.24 Miles	
Hungry after travelling so much? You can eat here:								
CATEGORY	NAME	ADDRESS	ZIPCODE	PHONE NO.	URL	ADDITIONAL D	Distance	Empi
Restaurant	CHICKEN, FR	PENN PLZ, M	10121	Not Available	Not Available	Not Available	0.16 Miles	
Restaurant	CHICKEN, CE	PENN PLAZA	10121	Not Available	Not Available	Not Available	0.16 Miles	
Restaurant	BEER PUB (E	PENN PLZ, M	10121	Not Available	Not Available	Not Available	0.16 Miles	
Restaurant	DUNKIN', HUI	PENNSYLVIA	10121	Not Available	Not Available	Not Available	0.16 Miles	
Restaurant	GARDEN MAI	PENN PLZ, M	10121	Not Available	Not Available	Not Available	0.16 Miles	
Restaurant	GARDEN MAI	PENN PLZ, M	10121	Not Available	Not Available	Not Available	0.16 Miles	
Restaurant	GARDEN MAI	PENN PLZ, M	10121	Not Available	Not Available	Not Available	0.16 Miles	
Restaurant	HOT DOG CO	PENN PLZ, M	10121	Not Available	Not Available	Not Available	0.16 Miles	
Restaurant	TACODUMBO	PENN PLAZA	10121	Not Available	Not Available	Not Available	0.16 Miles	
Restaurant	GARDEN MAI	PENN PLZ, M	10121	Not Available	Not Available	Not Available	0.16 Miles	

Visualizations –

The visualizations for our final project were achieved with the help of Tableau Software. The below visualization represents the concentration of the secondary tourist attractions close to the four primary tourist attractions in NYC (Empire State Building, New York Times Square, Central Park, and Statue of Liberty)



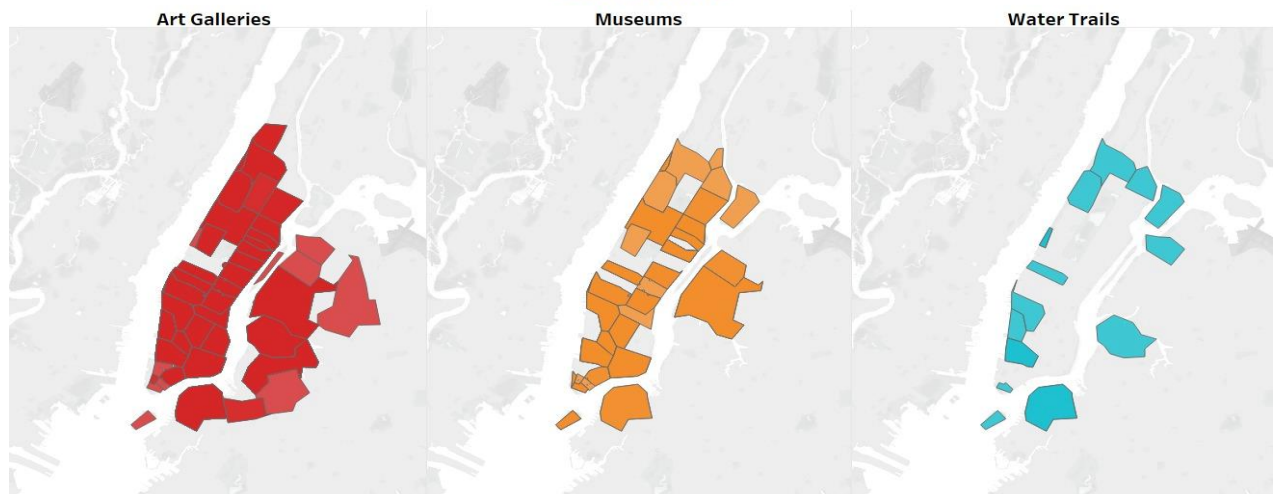
## Empire State Building



## Statue of Liberty



## Central Park



## **EXTERNAL ARTICLES**

Mentioned below are some of the external supporting materials that were used to provide additional insights into our business question and for our project:

- <https://www.planetware.com/tourist-attractions-/new-york-city-us-ny-nyc.htm>
- [https://www.tripadvisor.com/Attractions-g60763-Activities-New\\_York\\_City\\_New\\_York.html](https://www.tripadvisor.com/Attractions-g60763-Activities-New_York_City_New_York.html)
- <https://www.osc.state.ny.us/reports/osdc/tourism-industry-new-york-city>
- <https://www.tripsavvy.com/planning-a-new-york-city-trip-travel-guide-4178708>
- We could not find any other website or platform which performs a similar analysis. A trip planner called “Inspirock” (*Inspirock*) is quite close to what we are planning to achieve in this project where an itinerary is built for the user. However, it does not consider the type of tourists visiting New York City or the distance between different places.

Mentioned below are a few articles based on which we have allocated the tourist type for each of the secondary tourist locations:

- <https://www.timeout.com/new-york-kids/attractions/family-attractions>
- <https://secretnyc.co/things-to-do-with-friends-nyc/>
- <https://coupletraveltheworld.com/romantic-things-to-do-in-nyc-for-couples/>
- <https://www.theworldandthensome.com/15-fun-things-to-do-alone-in-nyc/>

Mentioned below is a link that helped us get the required zip codes for our project:

- [https://ny.postcodebase.com/zipcode\\_radius#myarticle](https://ny.postcodebase.com/zipcode_radius#myarticle)

## **ADDITIONAL DATA AND ANALYSIS**

- Our primary goal was to internally sort the restaurants based on the reviews fetched from the Yelp Datasets to recommend the best ones for the tourists. However, due to the huge file size, we faced technical difficulties in accessing the dataset using Python and converting it to CSV. Due to a shortage of time, we were unable to achieve this task.
- From the GUI perspective, we wanted to further customize and enhance the output to provide recommendations to the tourists in a more refined and attractive way. This can be considered one of the areas of improvement in the project.
- Post the completion of our project, we felt that including a few more secondary tourist attraction datasets around the four major ones could have resulted in a wider range of attraction spots for the tourists to choose from.