

Data Challenge

☰

Introduction

Quality Check

Data Munging

Data Visualization

» Key Factors

» Score Metric

» Map

Insights & What's Next

Introduction of Data Challenge

What is the problem?

Assumptions

Packages

The aim of the study is to assist a real estate company that has a niche in purchasing properties to rent out short-term as part of their business model specifically within New York City. The company has already figured out that 2-bedroom properties are the best for investment; however, they do not know which zipcodes are the best to invest in. A data product is requested to build in order to help the company understand which zipcodes would generate the most profit on short term rentals within New York City.

For this purpose, publicly available data from Zillow and AirBnB are used:

1. Revenue Data - AirBnB dataset with relevant short-term rental information  
<http://insideairbnb.com/get-the-data.html>
2. Cost Data - A historical estimated data of value for two-bedroom properties from 1996-04 to 2017-06, provided by Zillow  
<https://www.zillow.com/research/data/>
3. The Latest Cost Data - A historical estimated data of value for two-bedroom properties from 1996-04 to 2019-12, provided by Zillow  
<https://www.zillow.com/research/data/>

Data Challenge

☰

Introduction

Quality Check

Data Munging

Data Visualization

» Key Factors

» Score Metric

» Map

Insights & What's Next

Introduction of Data Challenge

What is the problem?

Assumptions

Packages

1. Short-term rental: A property that is rented anywhere between one evening up to one month is considered a short-term rental.
2. Weekly and monthly price: Weekly or monthly price of a property is a discount price for longer term stays. A property without a weekly or monthly price does not have discount price.
3. Booking habit: We assume that in a month 60% of bookings are single day bookings, 30% are weekly bookings and 10% of bookings are made for a month.
4. Occupancy rate: The occupancy rate is assumed to be 75%.
5. Weather/holiday has little or no impact on number of bookings.
6. The revenue obtained from a property remains equal to that charged by the previous host.
7. The company will put properties on rent throughout the year every day.
8. The investor will pay for the property in cash (i.e. no mortgage/interest rate will need to be accounted for).
9. The time value of money discount rate is 0% (i.e. \$1 today is worth the same 100 years from now).
10. All properties and all square feet within each locale can be assumed to be homogeneous (i.e. a 1000 square foot property in a locale such as Bronx or Manhattan generates twice the revenue and costs twice as much as any other 500 square foot property within that same locale.)

All the percentages assumed above can be modified by user preference.

Data Challenge

☰

Introduction

Quality Check

Data Munging

Data Visualization

» Key Factors

» Score Metric

» Map

Insights & What's Next

Introduction of Data Challenge

What is the problem?

Assumptions

Packages

Data preparation:

tidyverse, dplyr, stringr

Forecast:

forecast

Visualization:

RColorBrewer, plotly, leaflet

Dashboard:

shiny, shinydashboard

Data Challenge

☰

Introduction

Quality Check

Data Munging

Data Visualization

» Key Factors

» Score Metric

» Map

Insights & What's Next

Quality Check

Cost Data

Revenue Data

Get familiar with data

Data provides estimated historical median price for 2-bedroom homes in each zip code, captured between year 1996 and 2017 and spread monthly.

	RegionID	zipcode	1996-04	1996-05	1996-06	1996-07	1996-08	1996-09	1996-10	1996-11	1996-12	1997-01
1	61639	10025										
2	61637	10023										
3	61703	10128										
4	61625	10011										
5	61617	10003										

For future analysis, we filter the rows and select the relevant columns to keep the values of properties in New York City at the scraped date.

Data Quality

1. Missing values

Showing 1 to 6 of 6 entries

Previous1Next

For future analysis, we filter the rows and select the relevant columns to keep the values of properties in New York City at the scraped date.

### Data Quality

#### 1. Missing values

Median price for early years (1996-2013) has plenty of NAs as shown in the table below. Steps are taken in the following section to exclude columns with NAs when trend calculation and time series forecast.

Show10entries

Search:

	cost_missing_val
RegionID	0
zipcode	0
1996-04	17
1996-05	17
1996-06	17

#### 2. Missing zipcodes of NYC

New York hosts 176 zipcodes. However, there are only 25 zipcodes of NYC recorded in this data set. We are going to settle with this data for now. New data source can be connected in the future to account for rest of the zipcodes.

Data Challenge

Introduction

Quality Check

Data Munging

Data Visualization

>> Key Factors

>> Score Metric

>> Map

Insights & What's Next

## Quality Check

Cost Data

Revenue Data

### Get familiar with data

Revenue data contains information including details about the properties including location, number of bedrooms, room types, price and other details for stay.

Show10entries

Search:

id	listing_url	scrape_id	last_scraped	name	summary	space
----	-------------	-----------	--------------	------	---------	-------

In order to simplify the analysis, we only select 15 relevant variables and filter 2-bedroom properties in NYC to create a subset for the following section.

### Data Quality

#### 1. Needless character

### Data Quality

#### 1. Needless character

'\$' Value prefix of every price row prevents numeric manipulation. It is thus removed from three columns: Price, Weekly Price & Monthly Price.

#### 2. Missing values

Show10entries

Search:

	rev_missing_val
id	0
listing_url	0
scrape_id	0
last_scraped	0

#### (1) zipcode

There are less than 1% missing values in zipcode. Ignoring these values will not cause much effect to the analysis.

#### (2) weekly and monthly price

There are more than 80% weekly and monthly price missing. Since we assume that the weekly or monthly price of a property is a discount price for longer term stays, and a property without a weekly or monthly price does not have discount price, we can simply fill in the NAs with 7 and 31 times of daily price.

#### 3. Extreme values

There are a few extreme values in price such as \$9,000 or \$0 per night. The reason for this issue could be a data input error of the hosts. We can simply remove these rows.

#### 4. Ambiguous match between zipcode and neighborhood

#### 4. Ambiguous match between zipcode and neighborhood

Some zipcodes belong to more than one neighborhood in the dataset. To solve this problem, we keep the neighborhood to which the zipcode usually belongs as the only match.

## Data Challenge



Introduction

Quality Check

Data Munging

Data Visualization



» Key Factors

» Score Metric

» Map

Insights & What's Next

## Data Munging

Forecast

Data Join

Price Corrected

### Time Unit

The Zillow and Airbnb datasets have different units of time. In order to operate the analysis, we determine a common unit of time - month.

### ARIMA Forecast

The revenue data is scraped in 2019-07, but the cost data only have historical house value from 1996-04 to 2017-06. Due to shortage of time, we assume that there is seasonality in the price and that values depend not only on previous values (Auto Regressive AR) but also on differences between previous values (Moving Average MA). So, we apply, **ARIMA** model to predict the cost of the properties in zipcodes from 2017-07 to 2019-07. We then attach the price of property in 2019-07 calculated at zipcode level with each zipcode as a new column to cost data.

## Data Challenge



Introduction

Quality Check

Data Munging

Data Visualization



» Key Factors

» Score Metric

» Map

Insights & What's Next

## Data Munging

Forecast

Data Join

Price Corrected

In order to make sure this product is always applicable whenever new data is available or whenever we are ready to approach a new market, a function is built to link the data together in a scalable way by matching the scraped date and zipcode of two datasets.

After inner join, we recognize that there are 1508 properties of 24 zipcodes in 4 distinct neighborhoods - Manhattan, Queens, Brooklyn and Staten Island - matching together.

Data Challenge

Introduction
Quality Check
Data Munging
Data Visualization
Key Factors
Score Metric
Map
Insights & What's Next

## Data Munging

Forecast
Data Join
Price Corrected

### 1. Property type issue

If the property type == Private Room, the price of the property should be corrected by price\*bedrooms

### 2. Average daily price

In order to calculate the average daily price, we need to consider both tenant booking habit and host requirement for minimum and maximum nights.

Based on the common sense of general booking habit, we assume that in a month 60% of bookings are single day bookings, 30% are weekly bookings and 10% of bookings are made for a month.

As for the minimum/maximum amount of nights the host is willing to rent out the property, we assume that the minimum/maximum nights determine whether a property can be directly booked for a weekly and for a month.

Considering these two factors, we can calculate the average daily price by a conditional metric of minimum/maximum nights and booking habit weights.

Data Challenge

Introduction
Quality Check
Data Munging
Data Visualization
Key Factors
Score Metric
Map
Insights & What's Next

## Key factors

Zipcode	Neighbourhood	Number of Properties
11215	Manhattan	~180
10036	Manhattan	~145
10003	Manhattan	~135
11217	Manhattan	~125
10025	Manhattan	~115
10013	Manhattan	~105
10011	Manhattan	~100
11231	Brooklyn	~90
10014	Manhattan	~85
11201	Brooklyn	~85
10023	Manhattan	~75
10022	Manhattan	~70
10128	Manhattan	~65
10028	Manhattan	~55
10021	Manhattan	~35

### Controls

Number of zipcodes:

Key factors:

### Description

Higher the number of properties, more choices our client can invest, and more rental activities can possibly have in the area.

At the neighborhood level, Manhattan and Brooklyn host highest number of properties. Queens and Staten Island have much fewer number of properties.

Zipcodes 11215, 10036, 10003, 11217, 10025, 10013 have more than 100 properties in each area, which makes them the top 10 based on volume of properties.

Data Challenge

Introduction
Quality Check
Data Munging
Data Visualization
Key Factors
Score Metric
Map
Insights & What's Next

## Key factors

Zipcode	Neighbourhood	Cost of Properties
10013	Manhattan	~2.5M
10011	Manhattan	~2.2M
10014	Manhattan	~2.1M
10003	Manhattan	~1.8M
10028	Manhattan	~1.6M
10023	Manhattan	~1.5M
10021	Manhattan	~1.4M
10128	Manhattan	~1.3M
10022	Manhattan	~1.2M
10036	Manhattan	~1.1M
10025	Manhattan	~1.0M
11201	Brooklyn	~1.2M
11217	Brooklyn	~1.1M
11231	Brooklyn	~1.0M
11215	Brooklyn	~0.9M

### Controls

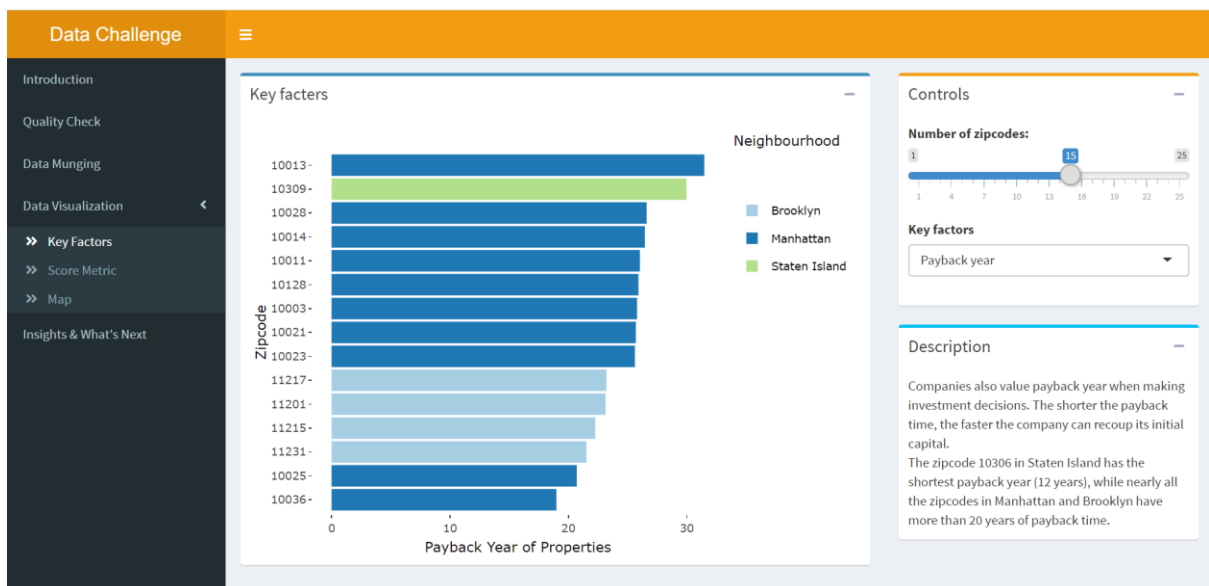
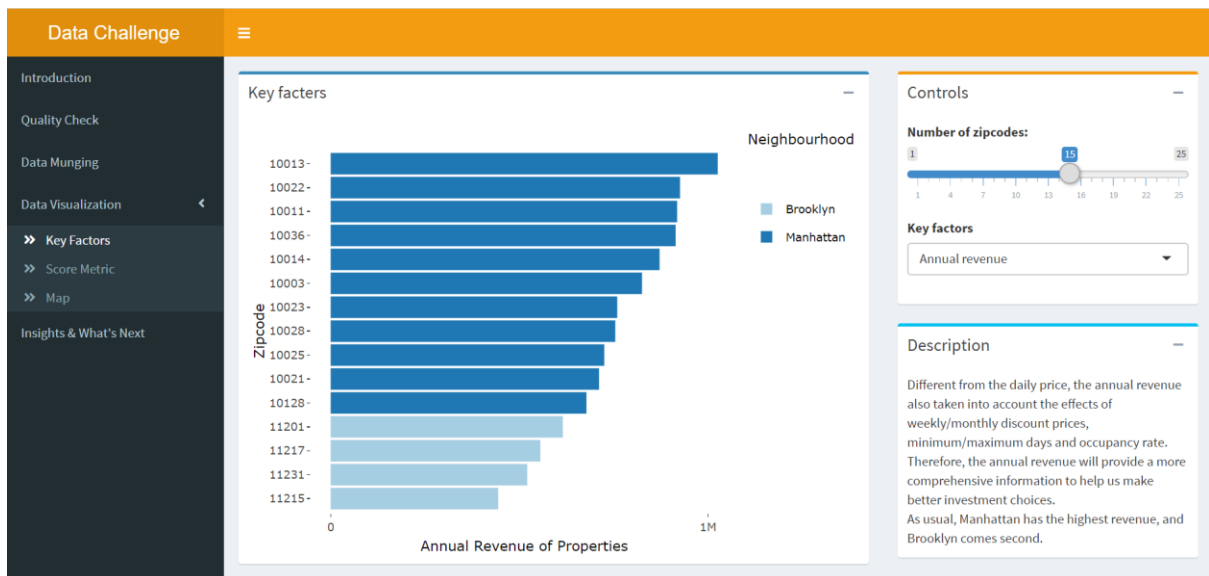
Number of zipcodes:

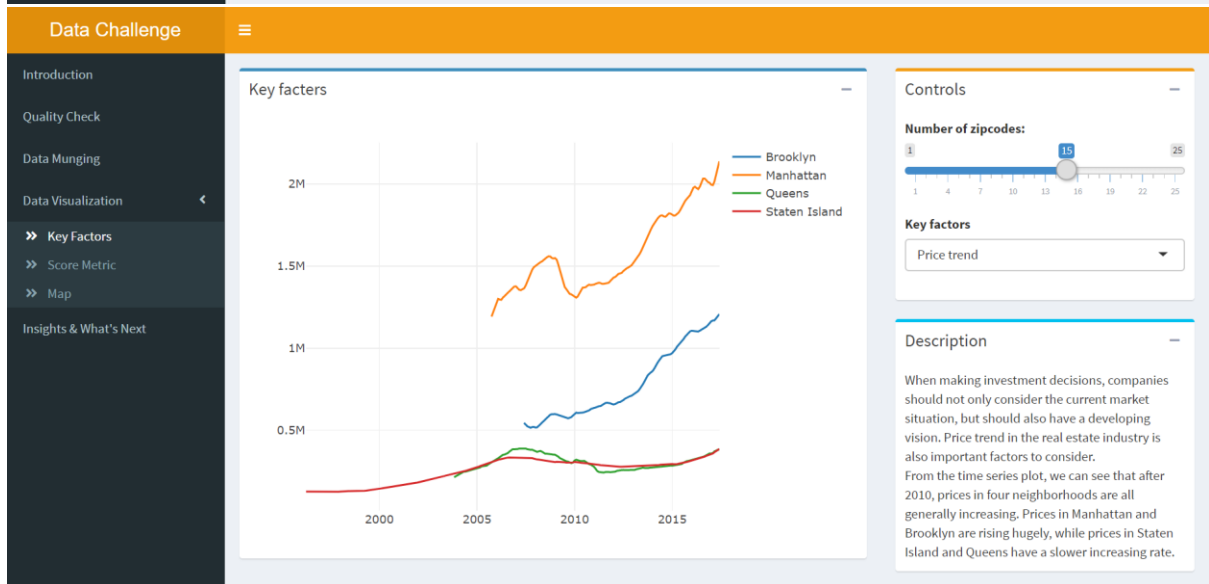
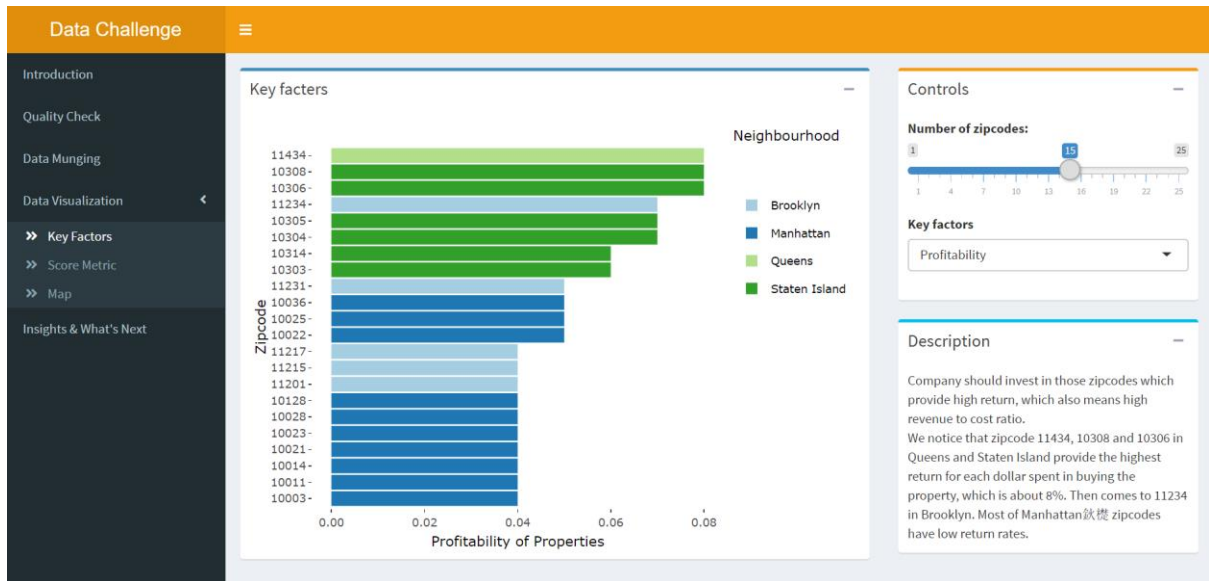
Key factors:

### Description

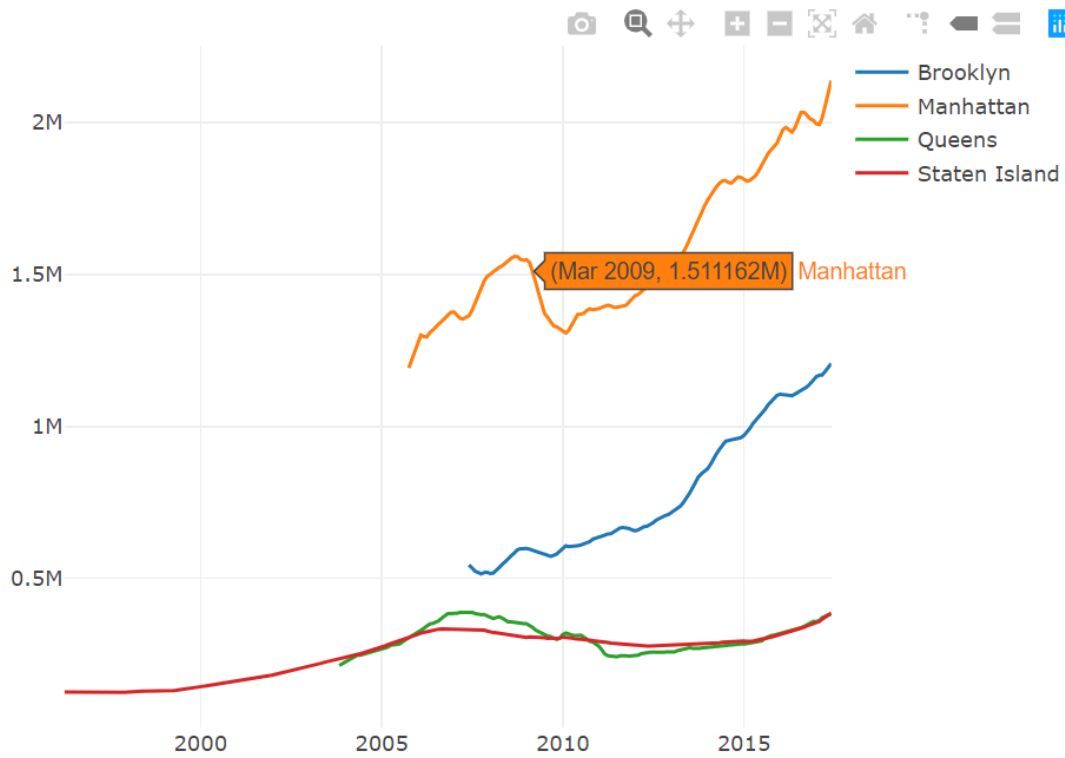
One of the primary constraints in decision of investment is the cost. Lower the cost, company can save a lot more and reach ROI faster.

Zipcodes in Manhattan walk away with highest property cost, with an average of about \$2M. Brooklyn comes second, with about \$1M. Zipcodes in Staten Island and Queens have much lower prices, with lower than \$0.5M.





## Key factors

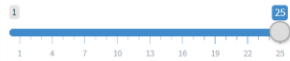


## Data Challenge

- Introduction
- Quality Check
- Data Munging
- Data Visualization
- Key Factors
- Score Metric
- Map
- Insights & What's Next

### Controls

Rows to show:



Columns to show:

- ☒ zipcode
- ☒ count
- ☐ avg\_monthly\_rev
- ☐ avg\_annual\_rev
- ☐ avg\_cost
- ☐ payback\_year
- ☐ rev\_cost\_ratio
- ☐ trend
- ☒ neighbourhood
- ☒ score

Weight of key factors

Count

Cost

score table score plot

Show 10 entries

Search:

	zipcode	count	neighbourhood	score
1	10003	134	Manhattan	1.15
2	10011	100	Manhattan	1.11
3	10013	103	Manhattan	1.13
4	10014	87	Manhattan	1.08
5	10021	27	Manhattan	0.76
6	10022	69	Manhattan	1.03
7	10023	77	Manhattan	0.95
8	10025	116	Manhattan	1.19
9	10028	54	Manhattan	0.86
10	10036	146	Manhattan	1.32

Showing 1 to 10 of 24 entries

Previous

1

2

3

Next



Insights & What's Next

☐ avg\_annual\_rev

☐ avg\_cost

☐ payback\_year

☐ rev\_cost\_ratio

☐ trend

☒ neighbourhood

☒ score

Weight of key factors

Count

0.2

Cost

0.2

Revenue

0.1

Profit

0.3

Payback year

0.1

Trend

0.1

4	10014	81	Manhattan	1.08
5	10021	27	Manhattan	0.76
6	10022	69	Manhattan	1.03
7	10023	77	Manhattan	0.95
8	10025	116	Manhattan	1.19
9	10028	54	Manhattan	0.86
10	10036	146	Manhattan	1.32

Showing 1 to 10 of 24 entries

Previous123Next

Description

Since there are multiple factors mentioned above influencing investment decisions, we can construct a scoring metrics to comprehensively evaluate these factors and score the different factors according to the set weights, and finally get a normalized score to select the best zipcodes.

When the weight of count, cost, revenue, profit, payback year and trend are initially set as 0.2, 0.2, 0.1, 0.3, 0.1, 0.1, the top 10 zipcodes with the highest score are: 11434, 10306, 11215, 10036, 10303, 10304, 10308, 10305, 10314, 11234. Six of them are located in Staten Island, two are in Brooklyn, one in Manhattan and one in Queens.

Data Challenge

Introduction

Quality Check

Data Munging

Data Visualization

Key Factors

Score Metric

Map

Insights & What's Next

Controls

Rows to show:

125

Columns to show:

☒ zipcode

☒ count

☐ avg\_monthly\_rev

☐ avg\_annual\_rev

☐ avg\_cost

☐ payback\_year

☐ rev\_cost\_ratio

☐ trend

☒ neighbourhood

☒ score

Weight of key factors

Count

Cost

score table

score plot

Neighbourhood

Brooklyn

Manhattan

Queens

Staten Island

11434-	
11215-	
10306-	
10036-	
10303-	
10308-	
10304-	
10305-	
11234-	
10314-	
11217-	
10025-	
11231-	
10003-	
10013-	
10011-	
11201-	
10014-	
10022-	
10023-	
10309-	
10128-	
10028-	
10021-	

Score of Properties

Data Challenge

Introduction

Quality Check

Data Munging

Data Visualization

Key Factors

Score Metric

Map

Insights & What's Next

Controls

Rows to show:

11325

Columns to show:

☒ zipcode

☒ count

☐ avg\_monthly\_rev

☐ avg\_annual\_rev

☐ avg\_cost

☐ payback\_year

☐ rev\_cost\_ratio

☐ trend

☒ neighbourhood

☒ score

Weight of key factors

Count

Cost

score table

score plot

Neighbourhood

Brooklyn

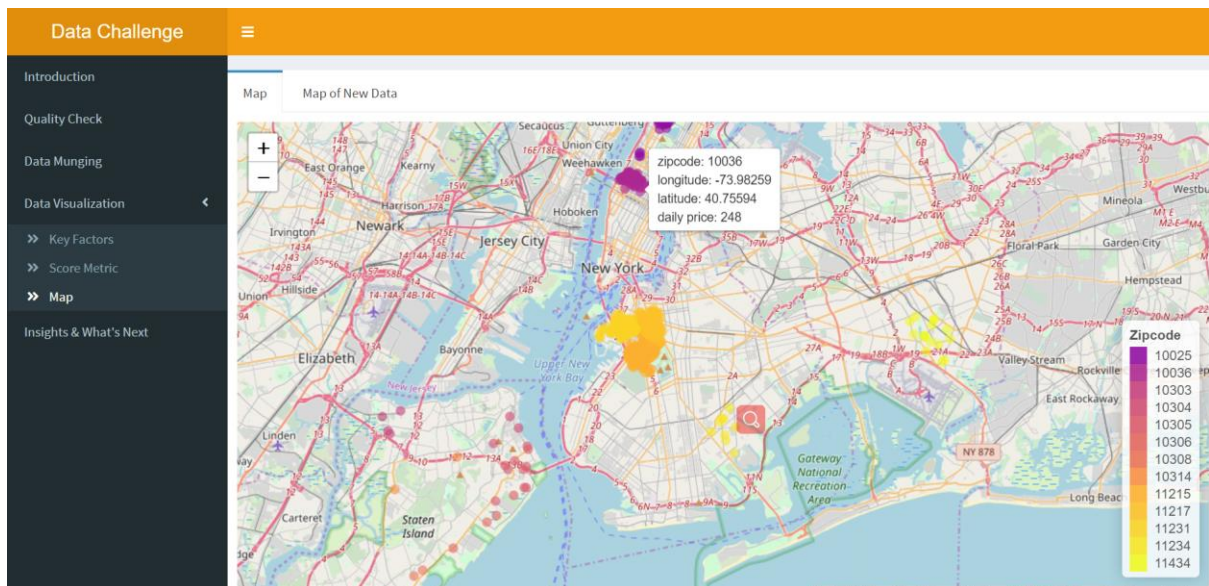
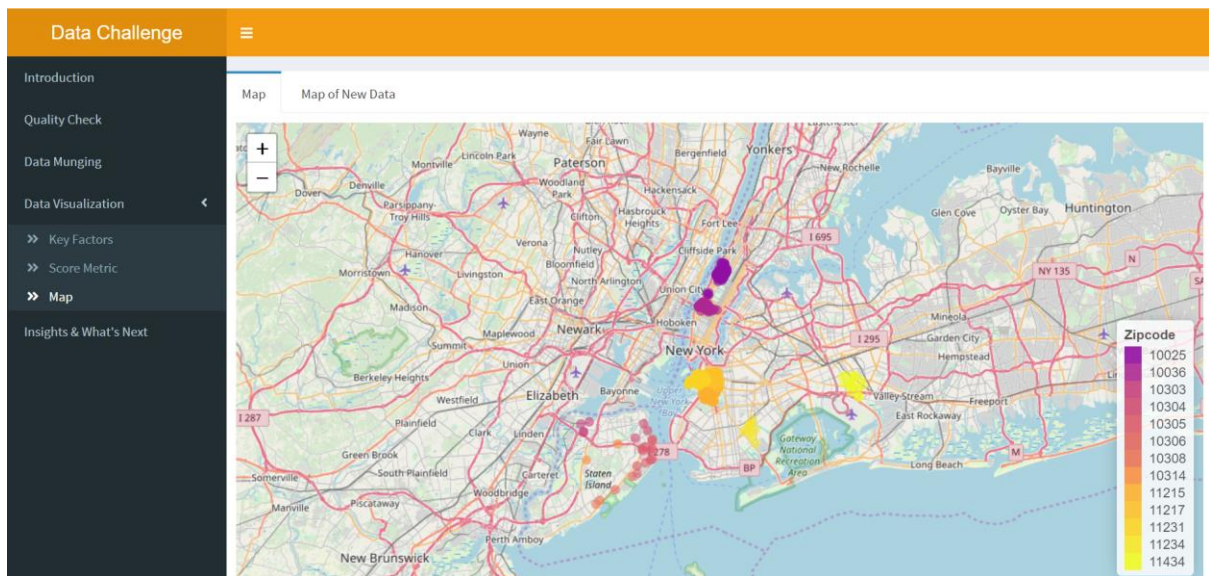
Manhattan

Queens

Staten Island

11434-	
11215-	
10306-	
10036-	
10303-	
10308-	
10304-	
10305-	
11234-	
10314-	
11217-	
10025-	
11231-	

Score of Properties



Data Challenge

Introduction

Quality Check

Data Munging

Data Visualization

» Key Factors

» Score Metric

» Map

Insights & What's Next

Map

Map of New Data

New data input

Browse...

No file selected

You can input 'Zip\_Zhvi\_2bedroom\_add.csv' in the folder. Wait a sec... Something magical is going to happen.

Data Challenge

Introduction

Quality Check

Data Munging

Data Visualization

Key Factors

Score Metric

Map

Insights & What's Next

Map

Map of New Data

New data input

Browse...

Zip\_Zhvi\_2bedroom\_add.csv

Upload complete

You can input 'Zip\_Zhvi\_2bedroom\_add.csv' in the folder. Wait a sec... Something magical is going to happen.

Data Visualization

Key Factors

Score Metric

Map

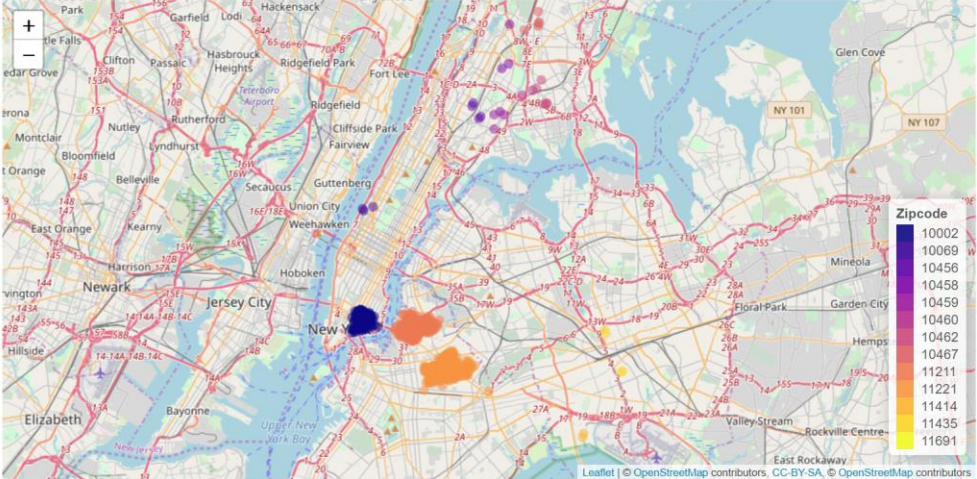
Insights & What's Next

Browse...

Zip\_Zhvi\_2bedroom\_add.csv

Upload complete

You can input 'Zip\_Zhvi\_2bedroom\_add.csv' in the folder. Wait a sec... Something magical is going to happen.



Data Challenge

Introduction

Quality Check

Data Munging

Data Visualization

Key Factors

Score Metric

Map

Insights & What's Next

Insights

What's Next

1. Comprehensive consideration:

Based on the comprehensive scoring method, the top 10 zipcodes with the highest score are: **11434, 10306, 11215, 10036, 10303, 10304, 10308, 10305, 10314, 11234**. Six of them are located in Staten Island, two are in Brooklyn, one in Manhattan and one in Queens.

However, after we add a latest dataset with more complete zipcodes, we can see a difference of the investment choices. The new top 10 zipcodes with the highest score are: **10069, 10459, 10460, 10462, 10467, 11211, 11221, 11414, 11435 and 11691**. Four of them are located in the Bronx, three are in Queens, two in Brooklyn and one in Manhattan. Interestingly, there is no property in Staten Island picked after inputting the new data.

This finding shows the importance of data integrity and timeliness. Therefore, we should take advantage of this automated app to keep the data up to date and check the latest best investment portfolio.

2. Different key factors consideration:

(The following insights come from the analysis of the old dataset.)

Number of properties & Cost of properties

Zipcodes 11215, 10303, 11434, 10304 and 10306 have substantial number of properties at low cost, which are located in Brooklyn, Staten Island and Queens. So, if the company has budget constraints then they should invest in buying properties in these zipcodes.

Weight of count = 0.5, weight of cost = 0.5

Annual revenue & Revenue cost ratio

## 2. Different key factors consideration:

(The following insights come from the analysis of the old dataset.)

### Number of properties & Cost of properties

Zipcodes 11215, 10303, 11434, 10304 and 10306 have substantial number of properties at low cost, which are located in Brooklyn, Staten Island and Queens. So, if the company has budget constraints then they should invest in buying properties in these zipcodes.

*Weight of count = 0.5, weight of cost = 0.5*

### Annual revenue & Revenue cost ratio

If the company is willing to buy properties having high cost, then they should invest in Zipcodes 10038 in Staten Island, 11215 in Brooklyn, 10036, 10003 and 10025 in Manhattan because these zipcodes not only have high number of costly properties but also provide very high return as the revenue is high.

*Weight of count = 0.5, weight of revenue = 0.2, weight of profit = 0.5*

### Payback year

If the company is more focused on getting its money back quickly with a limited initial investment, zipcodes 10038, 10314, 11434, 10303 and 10306 fit the bill, which are all located in Staten Island and Queens.

*Weight of cost = 0.5, weight of payback year = 0.5*

### Price trend

If the company attaches importance on long-term growth and future earnings and wants to invest in properties where rents are rising faster, these are the zipcodes to go: 10013, 10014, 10011 in Manhattan, 10308 in Staten Island, and 11201 in Brooklyn.

*Weight of profit = 0.5, weight of trend = 0.5*

## 3. Other thoughts:

- Our basic recommendation for the company is to diversify and buy properties in top performing zips of different neighborhoods with prime focus on Staten Island.
- Through our products, the company can adjust the weight of each factor according to its actual needs to obtain the most suitable investment portfolio plan for its own investment strategy.

## 3. Other thoughts:

- Our basic recommendation for the company is to diversify and buy properties in top performing zips of different neighborhoods with prime focus on Staten Island.
- Through our products, the company can adjust the weight of each factor according to its actual needs to obtain the most suitable investment portfolio plan for its own investment strategy.
- If the company is a risk taker, it can choose the property with high investment and high yield. If the company is risk averse, it can choose the property with low initial cost and fast return.
- Company should diversify its investment portfolios and locations to not only minimize investment risk, but also lay the foundation for future expansion in multiple markets.

## Data Challenge



Introduction

Quality Check

Data Munging

Data Visualization

» Key Factors

» Score Metric

» Map

Insights & What's Next

## Insights & What's Next

Insights

What's Next

- New York hosts 176 zipcodes, data can be further enriched to account for rest of the zipcodes. This would give the company more opportunities to diversify the investment portfolio.
- Future revenue prediction: In the present analysis, we used average daily price \* 360 \* occupancy rate to calculate the annual revenue. It is true based on the assumption that the rental revenue will not dramatically change within a year. However, since the house price is continuously increasing monthly, the rental price may also increase by time. Future analysis can consider the price trend of rental price in order to calculate income more accurately and realistically.
- Occupancy rate: In the present analysis, we simply assume the occupancy rate is 0.75 for every property in every zipcode. However, in reality, there may be many factors which can impact occupancy rate, such as review score, location and room quality. In the future, we could customize occupancy rate of each property by building a regression model.
- Interest rate: In this case, we have taken 0% discount rate as our assumption, but this assumption is not realistic. Some reasonable percentage rate can be taken to calculate NPV value and make a more accurate prediction.
- Text analytics on ignored description columns from revenue data: This would open insights about other metrics that drive customer to book an Airbnb property for rental such as access to public transport, parking space, etc.