# CISC 878

# PROJECT REPORT

## An exploratory data analysis and visualization of crime in Vancouver

**Zili Luo**

**20001744**

# Introduction

This report is cover my data analytic project, which did an exploratory data analysis and visualization of crime in Vancouver from 2003 to 2017. It will cover My dataset, Initial preprocessing of the dataset, Visualization results, potential correlations between reality and data points, predictive model, preprocessing for the model, and the predictive results.
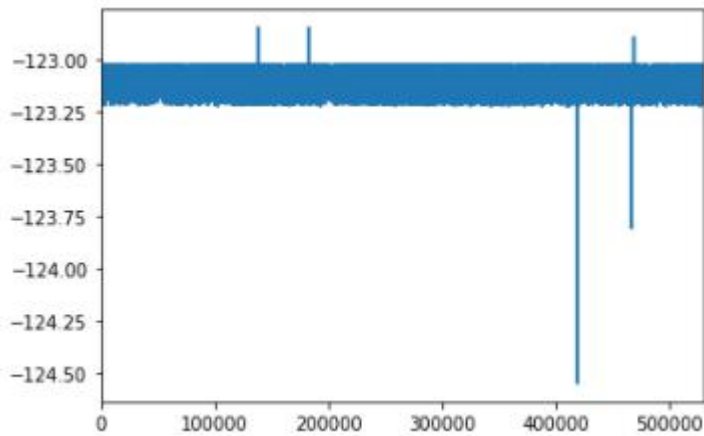
# Dataset

The dataset is available on Kaggle[1] named "crimes in Vancouver". It contains 530,652 records of crimes that happened in Vancouver between January 2003 and June 2017. Each row in the dataset includes the type of the crime, which has nine categories, when it has happened, the neighborhood where it has happened, and the spatial information represented in latitude and longitude. The dataset itself does not contain a rich information.In addition, about 10% of the data is anonymous, and 1% of the data is missing.

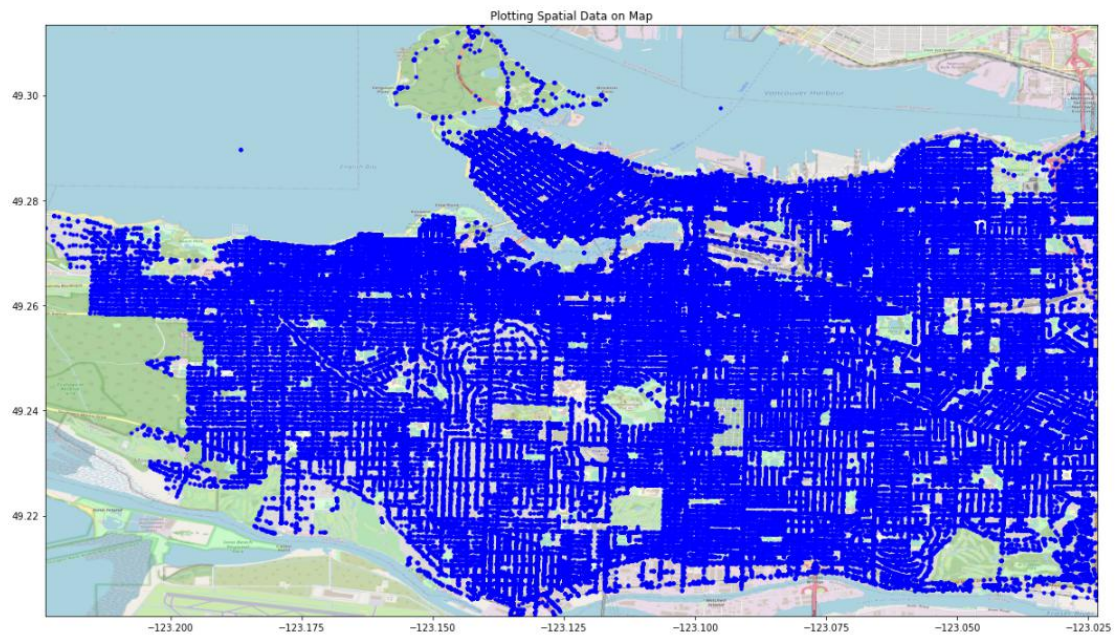| | TYPE | YEAR | MONTH | DAY | HOUR | MINUTE | HUNDRED_BLOCK | NEIGHBOURHOOD | X | Y | Latitude | Longitude |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Other Theft | 2003 | 5 | 12 | 16.0 | 15.0 | 9XX TERMINAL AVE | Strathcona | 493906.50 | 5457452.47 | 49.269802 | -123.083763 |
| 1 | Other Theft | 2003 | 5 | 7 | 15.0 | 20.0 | 9XX TERMINAL AVE | Strathcona | 493906.50 | 5457452.47 | 49.269802 | -123.083763 |
| 2 | Other Theft | 2003 | 4 | 23 | 16.0 | 40.0 | 9XX TERMINAL AVE | Strathcona | 493906.50 | 5457452.47 | 49.269802 | -123.083763 |
| 3 | Other Theft | 2003 | 4 | 20 | 11.0 | 15.0 | 9XX TERMINAL AVE | Strathcona | 493906.50 | 5457452.47 | 49.269802 | -123.083763 |
| 4 | Other Theft | 2003 | 4 | 12 | 17.0 | 45.0 | 9XX TERMINAL AVE | Strathcona | 493906.50 | 5457452.47 | 49.269802 | -123.083763 |
| 5 | Other Theft | 2003 | 3 | 26 | 20.0 | 45.0 | 9XX TERMINAL AVE | Strathcona | 493906.50 | 5457452.47 | 49.269802 | -123.083763 |
| 6 | Break and Enter Residential/Other | 2003 | 3 | 10 | 12.0 | 0.0 | 63XX WILTSHIRE ST | Kerrisdale | 489325.58 | 5452817.95 | 49.228051 | -123.146610 |

# Preprocessing

Since the data is simple, my preprocessing does not take many steps.
I dropped the missing data and anonymous data since both of them do not have a spatial information to visualize them. And I also dropped some outliers ,which are outside of Vancouver, I have found that might be caused by input mistakes.

# Visualization

The first visualization aims to find any hotspots that crimes are more likely to happen, which visualize all the data points on the map at once. But the number of data points is too large to find any trace; the result shows crimes happen everywhere in Vancouver during those years.



In the second visualization, the data is split by each year and also by the crime types. In this try, some traces of crime cluster could be observed in the map. All the map visualization result could be found at appendix.

Plotting Spatial Data of year 2003 on Map



Plotting Spatial Data of year 2016 on Map

The third visualization is on the data itself. Some histograms is created to find if anything in data itself are interesting in those grams.

Three histograms to plot crimes every year by a different type, hour, and neighborhood index.

crimes by hour

crimes by type

crimes by neighborhood

There isn't anything specific in the first two histograms, but in the crime by neighborhood gram, It seems in some neighborhoods, the crime rate has a significant change.

Then the total amount of crimes every year are plotted, since the data in 2017 is only a half-year data, the data in 2017 is doubled; we could observe that the number of crimes decreased from 2003 to 2010 and have gradually increased with a bump in 2016.

The next visualization is the theft crimes categories, which includes "Theft of Vehicle", "Other Theft", "Theft from Vehicle", "Theft of Bicycle", since the theft crimes occupies the majority of the dataset.

● First, lets start with "Theft of Vehicle":

Theft of Vehicle

- It had a major decrease, from an average of around 520 crimes per month in 2003 to around 100 in 2012.
- Although the average has been increasing in the past years, it's way below 2003.
- In 2002, the "Bait Car" program was launched and in 2003 the IMPACT group was formed in response to this peak in theft. It looks like they've been doing a great job!

- Second, about "Other Theft":



Other Theft

- On the opposite trend, other theft has been increasing, from around 200 to almost 500 crimes per month.
- With the decrease of "Theft of Vehicle", This means the theft still happens. While the vehicle itself might be well protected from being theft, the number of thieves does not decrease, and they still do their work.

- About "Theft from Vehicle":



Theft from Vehicle

- It is the most frequent type.
- It decreased along with "Theft of Vehicle" until 2012, but then it increased significantly.

- Finally, about "Theft of Bicycle":



- We can see a clear trend within the year. It has peaks during the summer months, which is expected.
- The average has also been increasing.

The last work I did in the visualization part is the Heat map. I have made two heat maps representing the average crimes that happen each day and average crimes in each hour.

The first heatmap to be introduced is average crimes in each day heatmap.

average crime heatmap

| DAY | Jan | Feb | Mar | Apr | May | Jun | Jul | Aug | Sep | Oct | Nov | Dec |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 93 | 93 | 92 | 97 | 93 | 92 | 99 | 99 | 103 | 94 | 103 | 91 |
| 2 | 85 | 93 | 80 | 90 | 88 | 90 | 91 | 94 | 92 | 92 | 89 | 82 |
| 3 | 89 | 89 | 81 | 91 | 84 | 88 | 92 | 91 | 95 | 91 | 94 | 90 |
| 4 | 84 | 90 | 85 | 90 | 84 | 90 | 95 | 94 | 90 | 92 | 87 | 88 |
| 5 | 87 | 89 | 88 | 87 | 87 | 93 | 95 | 92 | 97 | 93 | 96 | 89 |
| 6 | 87 | 89 | 86 | 87 | 96 | 94 | 90 | 97 | 94 | 87 | 95 | 90 |
| 7 | 98 | 90 | 89 | 91 | 93 | 97 | 92 | 96 | 95 | 97 | 96 | 90 |
| 8 | 94 | 86 | 88 | 96 | 98 | 95 | 89 | 100 | 98 | 89 | 94 | 84 |
| 9 | 91 | 86 | 83 | 97 | 93 | 92 | 93 | 99 | 92 | 92 | 90 | 85 |
| 10 | 90 | 86 | 87 | 89 | 96 | 94 | 86 | 100 | 100 | 92 | 97 | 87 |
| 11 | 86 | 83 | 91 | 98 | 93 | 97 | 92 | 99 | 101 | 92 | 100 | 90 |
| 12 | 85 | 90 | 104 | 92 | 92 | 99 | 97 | 98 | 104 | 94 | 97 | 85 |
| 13 | 89 | 92 | 93 | 98 | 94 | 98 | 92 | 94 | 97 | 89 | 91 | 89 |
| 14 | 88 | 90 | 88 | 95 | 100 | 94 | 92 | 101 | 95 | 97 | 100 | 85 |
| 15 | 96 | 96 | 91 | 91 | 102 | 132 | 97 | 103 | 98 | 104 | 104 | 85 |
| 16 | 91 | 84 | 88 | 89 | 93 | 97 | 98 | 96 | 100 | 99 | 98 | 83 |
| 17 | 94 | 84 | 93 | 91 | 95 | 99 | 91 | 97 | 101 | 95 | 88 | 76 |
| 18 | 83 | 84 | 89 | 90 | 93 | 100 | 94 | 94 | 98 | 92 | 86 | 76 |
| 19 | 83 | 86 | 95 | 90 | 88 | 90 | 94 | 94 | 91 | 87 | 86 | 78 |
| 20 | 83 | 82 | 93 | 81 | 97 | 102 | 95 | 98 | 95 | 93 | 89 | 80 |
| 21 | 78 | 76 | 86 | 83 | 92 | 93 | 93 | 100 | 92 | 91 | 82 | 71 |
| 22 | 77 | 79 | 89 | 81 | 90 | 84 | 95 | 98 | 90 | 91 | 81 | 72 |
| 23 | 81 | 75 | 81 | 78 | 88 | 89 | 86 | 96 | 89 | 88 | 79 | 82 |
| 24 | 78 | 80 | 81 | 80 | 91 | 86 | 88 | 88 | 84 | 86 | 75 | 80 |
| 25 | 78 | 76 | 76 | 77 | 79 | 78 | 92 | 86 | 83 | 87 | 80 | 57 |
| 26 | 81 | 78 | 83 | 81 | 78 | 83 | 84 | 87 | 85 | 84 | 76 | 73 |
| 27 | 84 | 79 | 80 | 76 | 84 | 84 | 80 | 79 | 85 | 80 | 81 | 82 |
| 28 | 83 | 80 | 81 | 79 | 81 | 87 | 82 | 84 | 83 | 88 | 83 | 82 |
| 29 | 86 | 86 | 86 | 80 | 82 | 84 | 81 | 88 | 85 | 88 | 80 | 86 |
| 30 | 87 | | 87 | 85 | 81 | 91 | 90 | 89 | 86 | 93 | 85 | 84 |
| 31 | 92 | | 89 | | 84 | | 86 | 88 | | 106 | | 87 |

Blue means good days. Red bad days. White average days.
The calmest day of crime is Christmas Day, December 25. That might because people all stay with their families, so no one walk in the street to be theft or criminals might also celebrate with families.
The worst day is June 15, which is caused by The Stanley Cup Riot in 2016. Over 600 crimes

happened on that day.

The first day of the month is usually a busy day for all months.

Halloween is also dangerous.

The second week of the summer months is usually the most dangerous.

BC Day (August 7) long weekend also has high averages.

The second heatmap to be introduced is average crimes in each hour heatmap.



crime hours heatmap

Most crimes happen between 5:00 pm and 8:00 pm

Break-and Enter-crimes happen almost all day.

All types of theft doesn't occur much in those between 3 am and 5 am as expected, that might because of most people are sleeping

After all the visualization, I tried to find if there are anything in reality could explain the changes in the dataset visualization.

At first, I checked the Vancouver Police Department website[4], seems they have the same strategy goal every year, which aims to decrease the amount of crime in Vancouver by 2.5%. But this can explain neither the increase from 2010 to 2016 and the bump in 2016. Then I searched for the big event that happened in Vancouver[3] and found Vancouver hosted the 2010 Winter Olympics and the election for the host city in 2003. This could explain the decrease from 2003 to 2010. And In August 2015, the 33-year-old businessman was killed by six Vancouver police officers. This news still has high clout in recent months. This event might cause a decrease in the trust of police at a significant amplitude.

Recall that the findings of some crimes are clustered in some neighborhoods. I went to the government of Vancouver website and found some evidence for this phenomenon[2]. The graph at right shows the average income of each area; we could observe the crime become sparser as the average income goes higher.



Also, except for the bump in 2016, the total crime curve each year has similar trends to the low-income rates curve in Vancouver[2].



## Potential relationship cont.

# predictive model

After all the visualization tasks, I'm curious about if a well-performed predictive model could be established to classify those crimes since there are some potential features of those criminal records. But as most data is related to theft, and the visualization results are not clear enough to support the idea. A bad result is expected.

The preprocessing part is still simple, as most features are already nominal. I merged some crime types, such as vehicle collision, break-and-enter. Select some essential elements such as month, hour, neighborhood, etc. And convert all ordinal data into categorical data.

```python
def category_merage(crime_type):
    if 'Break' in crime_type:
        return 'Break and Enter'
    elif 'Collision' in crime_type:
        return 'Vehicle Collision'
    else:
        return crime_type
```

The data feeding into the model is consist of hour when it has happened, the neighbourhood where it has happened, latitude and longtidude.

| DATE | HOUR | NEIGHBOURHOOD | Latitude | Longitude |
|------|------|---------------|----------|-----------|
| 2003-05-12 | 16.0 | 0 | 49.269802 | -123.083763 |
| 2003-05-07 | 15.0 | 0 | 49.269802 | -123.083763 |
| 2003-04-23 | 16.0 | 0 | 49.269802 | -123.083763 |
| 2003-04-20 | 11.0 | 0 | 49.269802 | -123.083763 |
| 2003-04-12 | 17.0 | 0 | 49.269802 | -123.083763 |
| 2003-03-26 | 20.0 | 0 | 49.269802 | -123.083763 |
| 2003-03-10 | 12.0 | 1 | 49.228051 | -123.146610 |
| 2003-06-28 | 4.0 | 2 | 49.255559 | -123.193725 |

The result of the predictive model is not good at all. With 38% training accuracy and similar testing accuracy. The model classified all the cases into class 1 and class 3, which are "Break and enter" and "Theft from vehicle". This result is under my expectation as those two cases happen everywhere and occupy the majority of the records.

```
              precision    recall  f1-score   support

           0       1.00      0.00      0.00     15480
           1       0.35      0.31      0.33     28669
           2       1.00      0.00      0.00     20844
           3       0.39      0.88      0.54     51428
           4       1.00      0.00      0.00      6703
           5       1.00      0.00      0.00     11339
           6       1.00      0.00      0.00      7742

    accuracy                           0.38    142205
   macro avg       0.82      0.17      0.12    142205
weighted avg       0.65      0.38      0.26    142205
```

I also modified the model several times but got a similar result. Sometimes the model can recognize the "Mischief" class, which also has a good proportion in the dataset and happens almost everywhere.

# Conclusion

The dataset itself does not contain that much information, but we could still observe some potential rules for crimes to happen via explore and visualization the data. And found something significant in reality have a correlation with crimes happening. Although it is nearly impossible to make a well-performed predictive model to predict when and where a crime will happen, which is also a problem in other crime data analytic with a dataset that contains hundreds of features, the police could still benefit from the findings to adjust their policy. That's what Intelligence-led policing encourages.

# Reference

[1] https://www.kaggle.com/wosaku/crime-in-vancouver

[2] https://vancouver.ca/files/cov/social-indicators-profile-city-of-vancouver.pdf

[3] https://dailyhive.com/vancouver/34-biggest-moments-vancouver-history

[4] https://vpd.ca/police/about/strategic-planning/index.html

# Appendix

Plotting Spatial Data of year 2005 on Map



Plotting Spatial Data of year 2006 on Map

Plotting Spatial Data of year 2007 on Map

Legend:
- Other Theft
- Break and Enter Residential/Other
- Mischief
- Break and Enter Commercial
- Theft from Vehicle
- Vehicle Collision or Pedestrian Struck (with Injury)
- Vehicle Collision or Pedestrian Struck (with Fatality)
- Theft of Vehicle
- Theft of Bicycle



Plotting Spatial Data of year 2008 on Map

Legend:
- Other Theft
- Break and Enter Residential/Other
- Mischief
- Break and Enter Commercial
- Theft from Vehicle
- Vehicle Collision or Pedestrian Struck (with Injury)
- Vehicle Collision or Pedestrian Struck (with Fatality)
- Theft of Vehicle
- Theft of Bicycle

Plotting Spatial Data of year 2009 on Map

Legend:
- Other Theft
- Break and Enter Residential/Other
- Mischief
- Break and Enter Commercial
- Theft from Vehicle
- Vehicle Collision or Pedestrian Struck (with Injury)
- Vehicle Collision or Pedestrian Struck (with Fatality)
- Theft of Vehicle
- Theft of Bicycle



Plotting Spatial Data of year 2010 on Map

Legend:
- Other Theft
- Break and Enter Residential/Other
- Mischief
- Break and Enter Commercial
- Theft from Vehicle
- Vehicle Collision or Pedestrian Struck (with Injury)
- Vehicle Collision or Pedestrian Struck (with Fatality)
- Theft of Vehicle
- Theft of Bicycle

Plotting Spatial Data of year 2011 on Map


Plotting Spatial Data of year 2012 on Map

Plotting Spatial Data of year 2013 on Map

Legend:
- Other Theft
- Break and Enter Residential/Other
- Mischief
- Break and Enter Commercial
- Theft from Vehicle
- Vehicle Collision or Pedestrian Struck (with Injury)
- Vehicle Collision or Pedestrian Struck (with Fatality)
- Theft of Vehicle
- Theft of Bicycle


Plotting Spatial Data of year 2014 on Map

Legend:
- Other Theft
- Break and Enter Residential/Other
- Mischief
- Break and Enter Commercial
- Theft from Vehicle
- Vehicle Collision or Pedestrian Struck (with Injury)
- Vehicle Collision or Pedestrian Struck (with Fatality)
- Theft of Vehicle
- Theft of Bicycle

Plotting Spatial Data of year 2015 on Map

Legend:
- Other Theft
- Break and Enter Residential/Other
- Mischief
- Break and Enter Commercial
- Theft from Vehicle
- Vehicle Collision or Pedestrian Struck (with Injury)
- Vehicle Collision or Pedestrian Struck (with Fatality)
- Theft of Vehicle
- Theft of Bicycle



Plotting Spatial Data of year 2016 on Map

Legend:
- Other Theft
- Break and Enter Residential/Other
- Mischief
- Break and Enter Commercial
- Theft from Vehicle
- Vehicle Collision or Pedestrian Struck (with Injury)
- Vehicle Collision or Pedestrian Struck (with Fatality)
- Theft of Vehicle
- Theft of Bicycle

Plotting Spatial Data of year 2017 on Map

Legend:
- Other Theft
- Break and Enter Residential/Other
- Mischief
- Break and Enter Commercial
- Theft from Vehicle
- Vehicle Collision or Pedestrian Struck (with Injury)
- Vehicle Collision or Pedestrian Struck (with Fatality)
- Theft of Vehicle
- Theft of Bicycle