

# Data

---

Based on definition of our problem, factors that will influence our decision are:

- finding the safest borough based on crime statistics
- finding the most common venues
- choosing the right neighbourhood within the borough

We will be using the geographical coordinates of Vancouver to plot neighbourhoods in a borough that is safe and in the city's vicinity, and finally cluster our neighborhoods and present our findings.

Following data sources will be needed to extract/generate the required information:

- **Part 1: Using a real world data set from Kaggle containing the Vancouver Crimes from 2003 to 2019**: A dataset consisting of the crime statistics of each Neighbourhood in Vancouver along with type of crime, recorded year, month and hour.
- **Part 2: Gathering additional information of the list of officially categorized boroughs in Vancouver from Wikipedia.**: Borough information will be used to map the existing data where each neighbourhood can be assigned with the right borough.
- **Part 3: Creating a new consolidated dataset of the Neighborhoods, along with their boroughs, crime data and the respective Neighbourhood's co-ordinates.**: This data will be fetched using OpenCage Geocoder to find the safest borough and explore the neighbourhood by plotting it on maps using Folium and perform exploratory data analysis.
- **Part 4: Creating a new consolidated dataset of the Neighborhoods, boroughs, and the most common venues and the respective Neighbourhood along with co-ordinates.**: This data will be fetched using Four Square API to explore the neighbourhood venues and to apply machine learning algorithm to cluster the neighbourhoods and present the findings by plotting it on maps using Folium.

Crimes in Vancouver Last Checkpoint: an hour ago (autosaved)

```
In [11]: vnc_crime_df['Neighbourhood'].value_counts()
Out[11]: Central Business District    10857
West End                            3031
Mount Pleasant                      2396
Strathcona                          1987
Kitsilano                           1802
Fairview                            1795
Renfrew-Collingwood                 1762
Grandview-Woodland                 1761
Kensington-Cedar Cottage            1391
Hastings-Sunrise                    1270
Sunset                              967
Riley Park                          866
Marpole                             828
Victoria-Fraserview                 600
Killarney                           565
Oakridge                            499
Dunbar-Southlands                   474
Kerrisdale                          417
Shaughnessy                         414
West Point Grey                     372
Arbutus Ridge                       311
South Cambie                        292
Stanley Park                        154
Musqueam                             17
Name: Neighbourhood, dtype: int64
```

```
In [5]: dfToronto = pd.DataFrame(index=Neighbourhoods, columns=["Population_2016", "Income_2016"])
dfToronto.head()
Out[5]:
```

	Population_2016	Income_2016
Aginicourt North	NaN	NaN
Aginicourt South-Malvern West	NaN	NaN

Crimes in Vanc...ipynb

Type here to search

Crimes in Vancouver Last Checkpoint: an hour ago (autosaved)

```
In [9]: vnc_crime_df = pd.read_csv('https://raw.githubusercontent.com/RamanujaSVL/Coursera_Capstone/master/vancouver_crime_records_2018.')
#Dropping X,Y which represents Lat, Lng data as Coordinates, the data seems to be corrupt
vnc_crime_df.drop(['Unnamed: 0', 'MINUTE', 'HUNDRED_BLOCK', 'X', 'Y'], axis = 1, inplace = True)
#vnc_crime_df.columns
vnc_crime_df.head()
Out[9]:
```

	TYPE	YEAR	MONTH	DAY	HOUR	NEIGHBOURHOOD
0	Break and Enter Commercial	2018	3	2	6	West End
1	Break and Enter Commercial	2018	6	16	18	West End
2	Break and Enter Commercial	2018	12	12	0	West End
3	Break and Enter Commercial	2018	4	9	6	Central Business District
4	Break and Enter Commercial	2018	10	2	18	Central Business District

```
In [10]: vnc_crime_df.columns = ['Type', 'Year', 'Month', 'Day', 'Hour', 'Neighbourhood']
vnc_crime_df.head()
Out[10]:
```

	Type	Year	Month	Day	Hour	Neighbourhood
0	Break and Enter Commercial	2018	3	2	6	West End
1	Break and Enter Commercial	2018	6	16	18	West End
2	Break and Enter Commercial	2018	12	12	0	West End
3	Break and Enter Commercial	2018	4	9	6	Central Business District
4	Break and Enter Commercial	2018	10	2	18	Central Business District

```
In [11]: vnc_crime_df['Neighbourhood'].value_counts()
Out[11]: Central Business District    10857
```

Crimes in Vanc...ipynb

Type here to search