# PREDICTING PRICES OF REAL-ESTATE IN CANADA

Data Science Project

By Omri Bakal

# RESEARCH QUESTION

- Can we predict the prices values of an apartment by other variables like : Number of beds,baths,name of the city?

# MY RESEARCH – APARTMENTS IN CANADA PROVINCES

I gathered information and data from the following provinces in canada:

Alberta

Newfoundland And Labrador

Ontario

British Columbia

Quebec

Prince Edward Island

Manitoba

Nova Scotia

New Brunswick

Sakatchewan

Data Source:

Crawling from: https://www.point2homes.com/CA

# MAIN STEPS

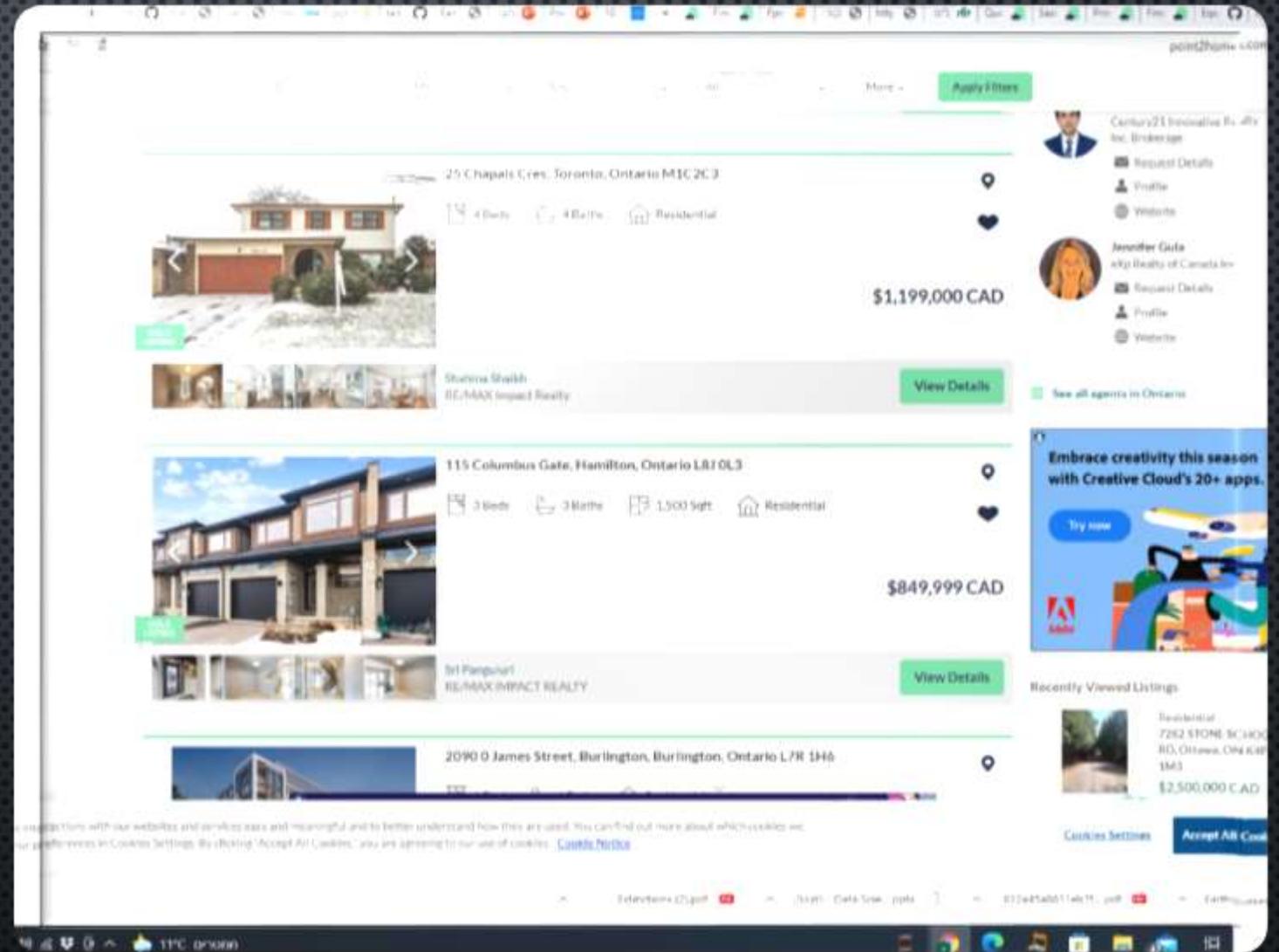Obtaining Data

Cleaning the Data

EDA Visualisation

Machine Learning

# CRAWLING FROM THE WEBSITE

I Took the most significant

Parameters from the website such as:

Number of beds,city,address,baths,

Seller,company,type of apartment and price.

Every Category of provence contains 30 pages of apartments , so I had to scrape 30 pages of all the provences.

PAGE OF ALL THE PROVENCES I HAD TO SCRAPE:

# THE FUNCTION THAT SCRAPING THE DATA AND PUTTING IT INTO A CSV:

- I PUT THE DATA INTO DIFFERENT LISTS BY PARAMETRS/CATEGORIES AND THEN TRANSFERRED IT INTO A CSV FILE.

# THE DATAFRAME:

I Gathered the data from 10 different categories of provences in Canada.

After collecting the data I transferred it into a csv file and then to a dataframe containing the columns:

- Address
- Price
- Type
- Beds
- Baths
- Company
- Seller
- city



```
In [495]: newdf=pd.read_csv("Apartments.csv")
          print(newdf)
          print(newdf.shape)
          print(newdf.isnull)
```

```
                                                    Address     Price  \
0              Colin Rd & Harmony Rd, Oshawa, Ontario L1K 1C1   1299.999
1                           6 Perfitt Cres, Ajax, Ontario L1Z1J5    899.900
2       Hwy 7 & Jane St, Vaughan, Vaughan, Ontario L4K...    499.900
3       Linea Condos   /743 Warden Ave, Scarborough, To...    569.900
4                556389 MULMUR/MEL TL LINE, Mulmur, Ontario   2988.000
...                                                     ...        ...
4609    46 Kirk CRESCENT, Saskatoon, Saskatchewan S7H 3B2    569.000
4610    601 1st AVENUE W, Zenon Park, Saskatchewan S0E...     96.000
4611    910 9th STREET E 106, Saskatoon, Saskatchewan ...    209.900
4612    25 Wellington DRIVE, Moose Jaw, Saskatchewan S...    269.900
4613    130 Marlatte CRESCENT 1204, Saskatoon, Saskatc...    225.900

                 Type  Beds  Baths  \
0         Residential     4      4
1         Residential     3      4
2         Residential     1      1
3         Residential     1      1
4       Single Family     5      3
```

# DATA CLEANING DATA TYPE CONVERTING

- As you can see in the following picture, I dropped all the missing values in the columns Beds and Baths , and I also converted the object categories Beds, Baths and Price into a int and float values.

- This step is going to help futher on EDA step and Machine Learning.

```
newdf=newdf.dropna(subset=["Beds","Baths"],axis=0)
newdf['Baths'] = newdf['Baths'].astype(int)
newdf['Beds'] = newdf['Beds'].astype(int)
newdf['Price']=newdf['Price'].astype(float)
newdf['Price']=newdf['Price']/1000
print(newdf)
```

```
                                            Address    Price  \
0        McLevin Ave and Tapscott Rd, Toronto, Toronto,...  799900
1        Sixth Line & Bowbeer Road, Oakville, Oakville,...  2199900
2        540 Davis Dr W, Newmarket, Newmarket, Ontario ...   550000
3        Cottonwood Cres Welland Ontario, Welland, Ontario  999900
4                       182 Verdun Rd, Oshawa, Ontario L1H5T2   499900
...                                                   ...      ...
6714   25 Wellington DRIVE, Moose Jaw, Saskatchewan S...   269900
```

# AS YOU CAN SEE IN THE DF.INFO() COMMAND, NOW THE DATA CONTAINS OBJECT COLUMNS AND ALSO INT/FLOAT COLUMNS

```
newdf.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 4664 entries, 0 to 6718
Data columns (total 8 columns):
 #   Column    Non-Null Count   Dtype
---  ------    --------------   -----
 0   Address   4664 non-null    object
 1   Price     4664 non-null    float64
 2   Type      4664 non-null    object
 3   Beds      4664 non-null    int32
 4   Baths     4664 non-null    int32
 5   Company   527 non-null     object
 6   Seller    528 non-null     object
 7   City      3365 non-null    object
dtypes: float64(1), int32(2), object(5)
memory usage: 291.5+ KB
```

# EDA VISUALISATIONS

Type of Apartment by Price

# TYPE OF APARTMENT BY PRICE GRAPH

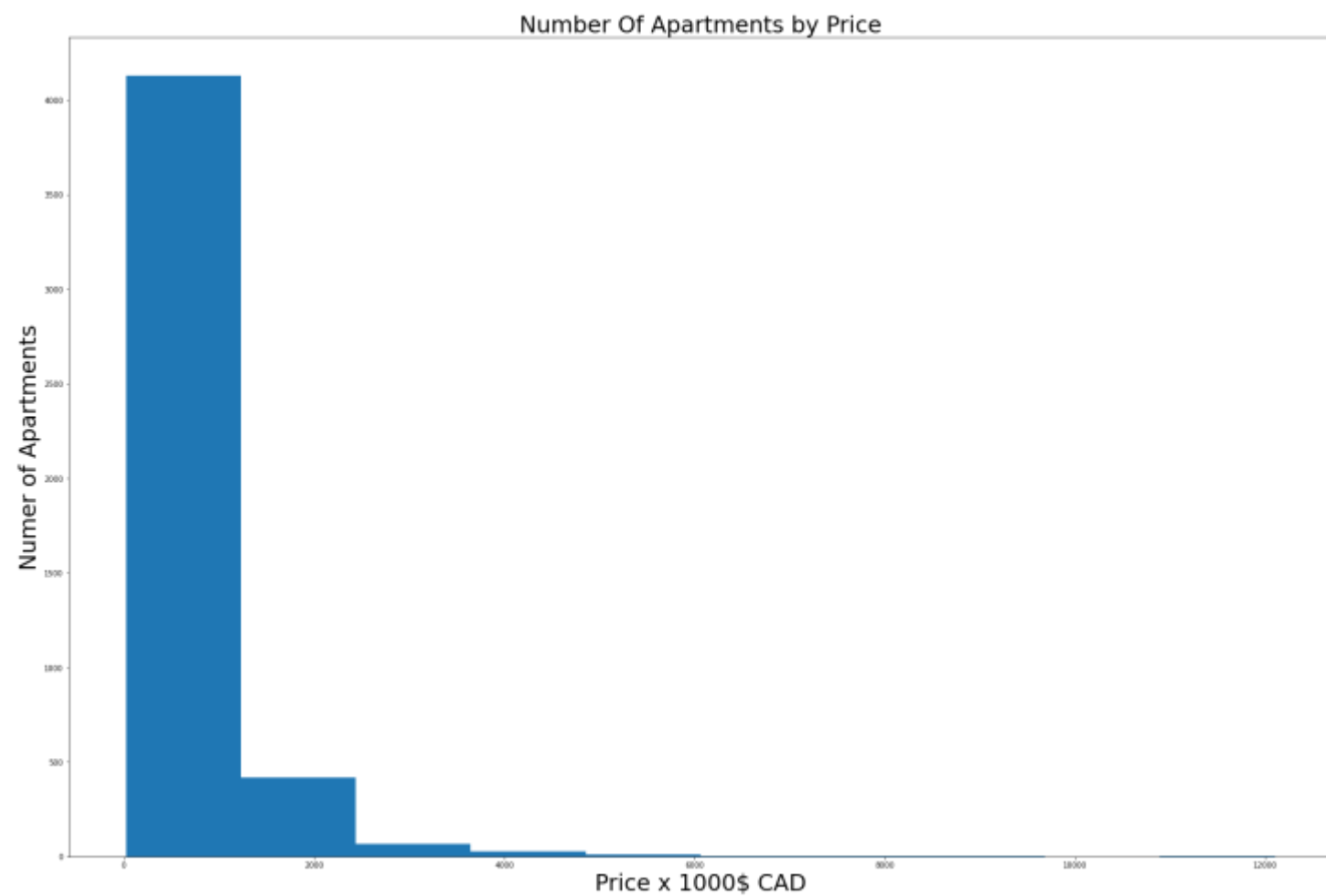- THIS IS A **BAR** PLOT , COMPARES THE AMOUNT TYPE OF HOUSES BY THEIR PRICES.

- X BAR CONTAINS THE TYPE OF THE RESIDENCE , AND THE Y BAR CONTAINS THE PRICE

- WE CAN SEE FROM THE PLOT THAR RESIDENTIAL TYPE OF HOUSES ARE MORE WANTED AMONG THE PEOPLE IN CANADA.

Number Of Apartments by Price
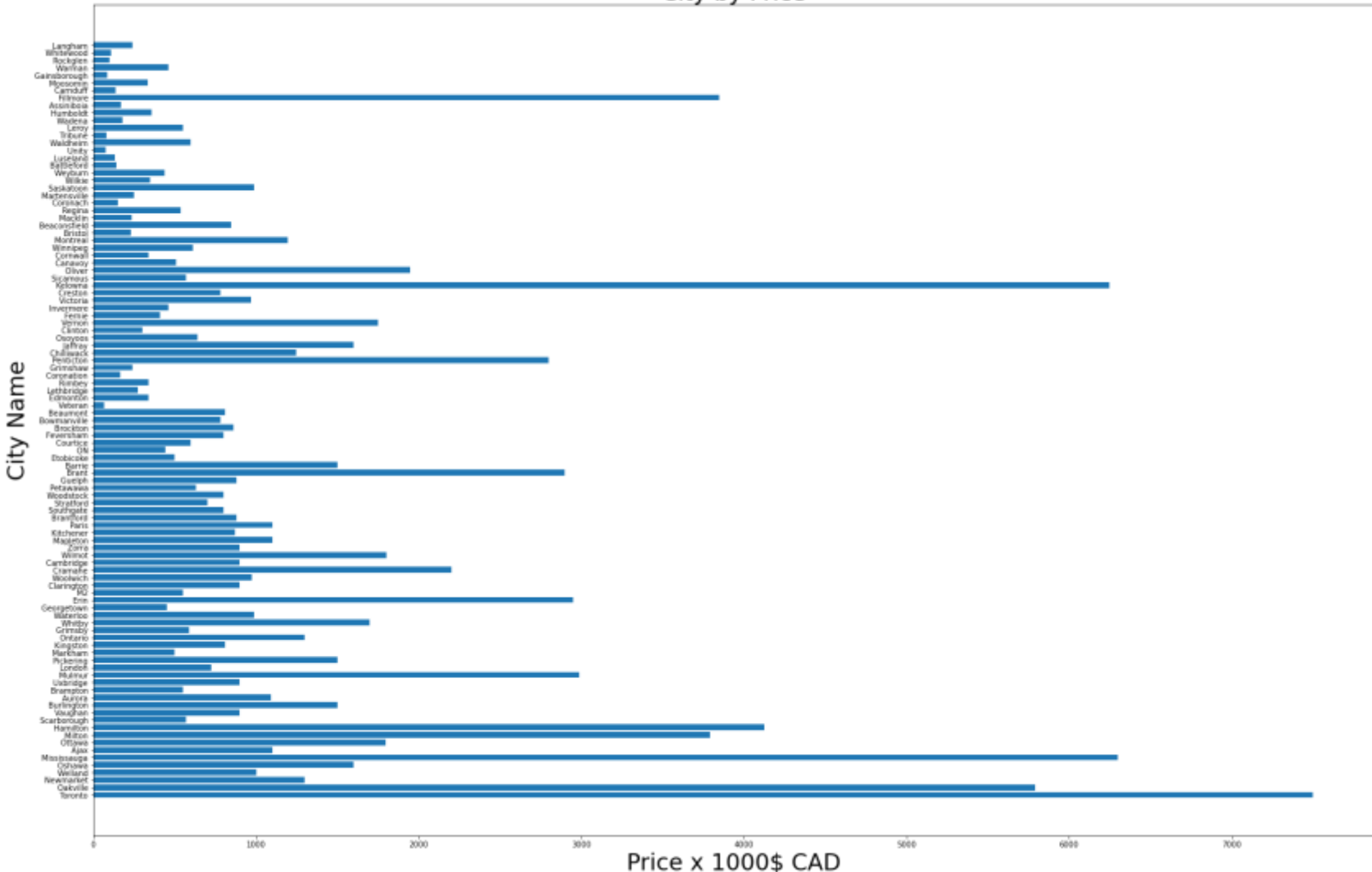
# NUMBER OF APARTMENTS BY PRICE GRAPH:

This is **Histogram** Plot , that shows us the number of the Apartments in Canada by their prices.

The X bar contians the Prices , and the Y bar contains the amount of the houses in the range of the price shown below.

We can see from this plot that most of the houses in the provinces in Canada are sold between 1-2 million CAD dollars.
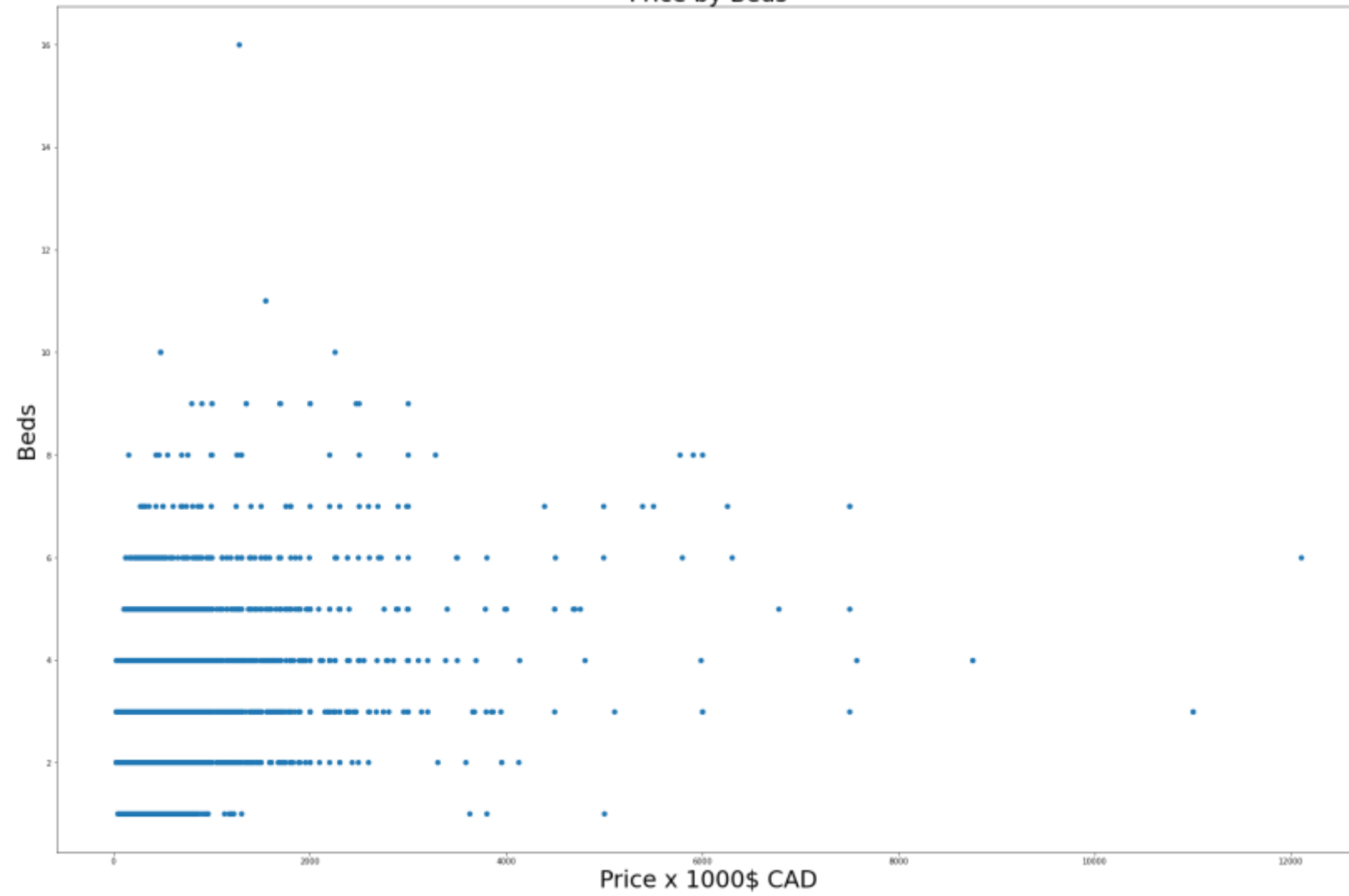
City by Price

# CITY BY PRICE GRAPH:

- This is another **BAR** plot that shows us which city in canada is the most expensive , we can see from the plot that toronto has the most expensive houses from all other cities.

- The x bar represents the price and the y bar represnts the city.

Price by Beds

# PRICE BY BEDS SCATTER PLOT:

- THIS IS A **SCATTER** EHT SU SWOHS TAHT TOLP EHT OT GNIDROCCA SECIRP EHT FO SEULAV GNIRETTACS .SDEB FO REBMUN

- THE X BAR REPRESENTS THE PRICES , AND THE Y BAR REPRESENTS THE PRICES.

- THIS IS AN IMPORTANT PLOT THAT WILL HELP US FUTHER ON IN THE MACHINE LEARNING PROCESS,WE CAN SEE FROM THIS PLOT THAT THERE IS A CORRELATION BETWEEN THE TWO VARIABLES —THE PRICE OF THE APARTMENT IS GETTING MORE EXPENSIVE BY THE NUMBER OF BEDS.

# MACHINE LEARNING

# MAIN MACHINE LEARNING TYPE:SUPERVISED LEARNING

Because the prices values are continious varibales , I had to use a supervised learning method to get the best predicted values , and I had to use two models during the process:

1.Linear Regression

2.Random Forest Regressor

## LINEAR REGRESSION:

- As you can see in the code below , I removed the problematic string columns from the DataFrame , and then splitted the DataFrame into X,y in a ratio of 80% and 20%.

- Unfortunately, the pridected model score had a low value so I had to change to other one.

```
In [362]: y=newdf.iloc[:,1].values
          X=newdf.drop(columns=["Address","Company","Seller","Price","Type","City"])
          X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=0)
          model=LinearRegression(fit_intercept=False)
          model.fit(X_train,y_train)
          y_pred=model.predict(X_test)
          print(r2_score(y_test,y_pred))

          0.30216115482533346
```
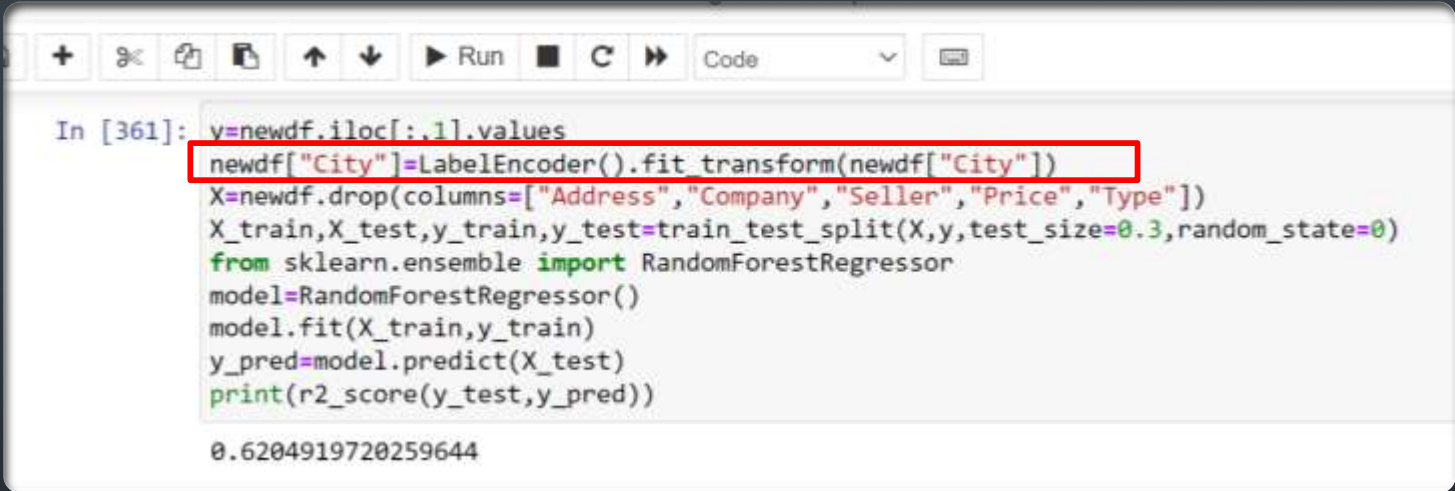
# RANDOM FOREST REGRESSOR:

So after getting low results on the linear regression model , I tried another supervised model and had to choose a random forest regressor.

In addition I used the Label Encoder feature of sklearn to transform the city coulmn into a numeric column , something that can help in predicting the model more precisely.

Seems like the score of the model got a better result after that.



```python
In [361]:  y=newdf.iloc[:,1].values
           newdf["City"]=LabelEncoder().fit_transform(newdf["City"])
           X=newdf.drop(columns=["Address","Company","Seller","Price","Type"])
           X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.3,random_state=0)
           from sklearn.ensemble import RandomForestRegressor
           model=RandomForestRegressor()
           model.fit(X_train,y_train)
           y_pred=model.predict(X_test)
           print(r2_score(y_test,y_pred))

           0.6204919720259644
```

# CONCLUSIONS

# FINAL WORDS...

SEEMS LIKE THE MODEL SUCCEEDED TO PREDICT 62% OF THE VALUES , ACCORDING TO THE VALUES OF NUMBER OF BEDS AND BATHS AND ALSO THE CITY.