

Team Members

- Alexis Lawal
- Dayo Thompson
- Kirushan Kirubaharan
- Sushant Deshpande

Objective

- Extracting Real Estate data in Calgary,
 Alberta and Walk Scores for corresponding house addresses
- Walk Score: shows a walk/transit/bike score (0-100) for any address to the local downtown
- Transforming retrieved data into easy-to-read tables
- Loading transformed data into relational and non-relational databases for optimal functionality

Data Sources

- Remax Canada: https://www.remax.ca/ab/calgary-real-estate
- Walk Score: https://www.walkscore.com/CA-AB/Calgary



Scraping Calgary Real Estate Data

Remax

fppt.com

```
In [1]: import pandas as pd
   import numpy as np
   import requests
   from bs4 import BeautifulSoup
   import time
   from splinter import Browser
   from sqlalchemy import create_engine
   import warnings
   warnings.filterwarnings('ignore')
   print('Libraries imported!')
```

Using BeautifulSoup to scrape property details (house address, house details).

```
In [2]: house address = []
        house details = []
        base url = 'https://www.remax.ca/ab/calgary-real-estate?page='
        urls = [base url + str(x) for x in range(1,301)]
        time.sleep(2)
        for url in urls:
            # Parse HTML with Beautiful Soup
            time.sleep(2)
            response = requests.get(url)
            soup = BeautifulSoup(response.text, 'html.parser')
            try:
                addresses = soup.find all('div', class = 'left-content flex-one')
                for address in addresses:
                    house address.append(address.text)
            except:
                house address.append('None')
            try:
                details = soup.find all('div', class = 'property-details')
                for detail in details:
                    house details.append(detail.text)
            except:
                house details.append('None')
```

Scraping Walk Score Data

```
In [ ]: scores walk = []
        scores bike = []
        scores transit = []
        for i in post code list:
            try:
                postal_code = i.replace(" ", "%20")
                url score = "https://www.walkscore.com/score/" + str(postal code)
                time.sleep(2)
                # Parse HTML with Beautiful Soup
                response = requests.get(url score)
                code soup = BeautifulSoup(response.text, 'html.parser')
                if 'pp.walk.sc/badge/walk/score' in str(code soup):
                    ws = str(code soup).split('pp.walk.sc/badge/walk/score/')[1][:2].replace('.','')
                    scores walk.append(ws)
                else:
                    ws = 'N/A'
                    scores walk.append(ws)
                if 'pp.walk.sc/badge/bike/score' in str(code soup):
                    bs = str(code soup).split('pp.walk.sc/badge/bike/score/')[1][:2].replace('.','')
                    scores bike.append(bs)
                else:
                    bs = 'N/A'
                    scores bike.append(bs)
                if 'pp.walk.sc/badge/transit/score' in str(code soup):
                    ts = str(code soup).split('pp.walk.sc/badge/transit/score/')[1][:2].replace('.','')
                    scores transit.append(ts)
                else:
                    ts = 'N/A'
                    scores transit.append(ts)
            except:
                ws = 'N/A'
                scores walk.append(ws)
                bs = 'N/A'
                scores bike.append(bs)
                ts = 'N/A'
                scores transit.append(ts)
```



Cleaning the Calgary Real Estate Data

First dataframe: Address and Price details

```
In [3]: address df = pd.DataFrame(house address)
         new df = address df[0].str.split(' ', 2, expand=True)
         new df["price"] = new df[1].str.replace("$", "")
         new df["price"] = new df["price"].str.replace(",", "")
         new df["price"] = pd.to numeric(new df["price"])
         del new df[0]
         del new df[1]
         new df.head()
Out[3]:
                                                     price
                 9803 ELBOW DR SW, Calgary, AB, T2V 1M4
                                                    489900
                 101 - 3704 15A ST SW, Calgary, AB, T2T 4C3
                                                    319900
          2 25 HARVEST GLEN WAY NE, Calgary, AB, T3K 4J2
                                                    399900
              32 EVERGLEN GROVE SW, Calgary, AB, T2Y 4Z3 429500
              416 THORNDALE RD NW, Calgary, AB, T2K 3C5 484900
```

fppt.com

Cleaning the Calgary Real Estate Data

- Values separated into columns: price, address, postal code, bedrooms, bath, property type
- Price column type changed to integer

```
In [5]: final df = new df[2].str.split(', Calgary, AB, ', expand=True)
         final df.head()
Out[5]:
                                         1
                  9803 ELBOW DR SW T2V 1M4
                 101 - 3704 15A ST SW T2T 4C3
            25 HARVEST GLEN WAY NE
                                  T3K 4J2
              32 EVERGLEN GROVE SW T2Y 4Z3
               416 THORNDALE RD NW T2K 3C5
         df add = pd.concat([new df, final df], axis=1)
In [6]:
         del df add(2)
         df add.columns = ["price", "address", "postal code"]
         df add.head()
Out[6]:
                                   address
                                          postal code
              price
            489900
                         9803 ELBOW DR SW
                                              T2V 1M4
            319900
                        101 - 3704 15A ST SW
                                              T2T 4C3
            399900
                    25 HARVEST GLEN WAY NE
                                              T3K 4J2
             429500
                     32 EVERGLEN GROVE SW
                                              T2Y 4Z3
                                              T2K 3C5
             484900
                      416 THORNDALE RD NW
```

Cleaning the Calgary Real Estate Data

Second dataframe: House details

```
In [6]: details = pd.DataFrame(house details)
         details df temp = details[0].str.split('|', expand=True)
         details df temp.head()
Out[6]:
                     2 bath
          0 4 bed
                           1121 sqft house
                     1 bath 836 sqft condo
            2 bed
          2 1 bed 1 bath 969 sqft house
         3 4 bed 2 bath 1650 sqft house
         4 4 bed 3 + 1 bath 2805 sqft house
```

Joining the Calgary Real Estate Dataframes

 Concatenating House Address/Price details and House details dataframes.

```
In [12]:
          calgary df dup = pd.concat([df add, details df], axis=1)
          calgary df = calgary df dup.drop duplicates()
          calgary df.head()
Out[12]:
                                                   postal code
                                                               bed_full_bath_half_bath_property_area_property_type
                price
           0 489900
                                9803 ELBOW DR SW
                                                      T2V 1M4
                                                                4.0
                                                                         2.0
                                                                                  NaN
                                                                                              1121 0
                                                                                                            house
              239900
                      106 - 790 KINGSMERE CRES SW
                                                      T2V 2G9
                                                                2.0
                                                                         1.0
                                                                                  NaN
                                                                                               836.0
                                                                                                            condo
           2 789900
                                                       T2T 4L7
                                                                         1.0
                                    4508 16A ST SW
                                                                1.0
                                                                                  NaN
                                                                                               969.0
                                                                                                            house
                                    1015 19 AVE SE
                                                      T2G 1M1
                                                                         2.0
           3 595000
                                                                40
                                                                                  NaN
                                                                                              1650 0
                                                                                                            house
             799900
                            96 ASPEN STONE RD SW
                                                      T3H 5Y7
                                                                4.0
                                                                         3.0
                                                                                   1.0
                                                                                              2805.0
                                                                                                            house
In [13]:
          calgary df.to csv('calgary df.csv', index=False)
```

Cleaning the Walk Score Data

 Data converted into dataframe and columns named.

In [9]:	score_df.head()						
Out[9]:		postal_code	walk_score	bike_score	transit_score		
19 (10 (19))	0	T2V 1M4	58.0	61.0	55.0		
	1	T2T 4C3	53.0	81.0	42.0		
	2	T3K 4J2	19.0	59.0	38.0		
	3	T2Y 4Z3	6.0	30.0	31.0		
	4	T2K 3C5	61.0	81.0	53.0		



Loading Data to Relational Database (PostgreSQL)

Building connection to PostgreSQL and loading transformed data

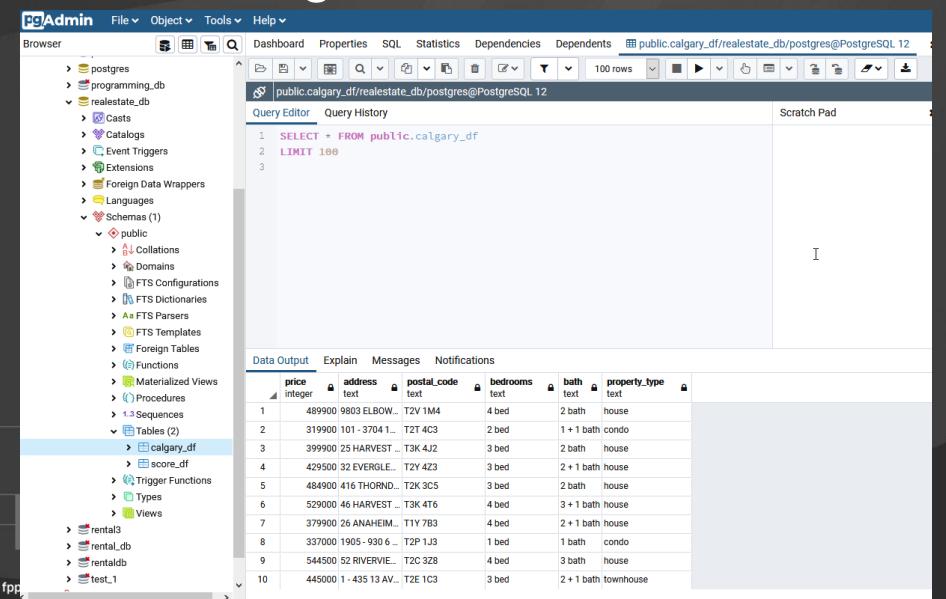
SQL

```
In [29]: calgary_df = pd.read_csv('calgary_df.csv')
    score_df = pd.read_csv('score_df.csv')

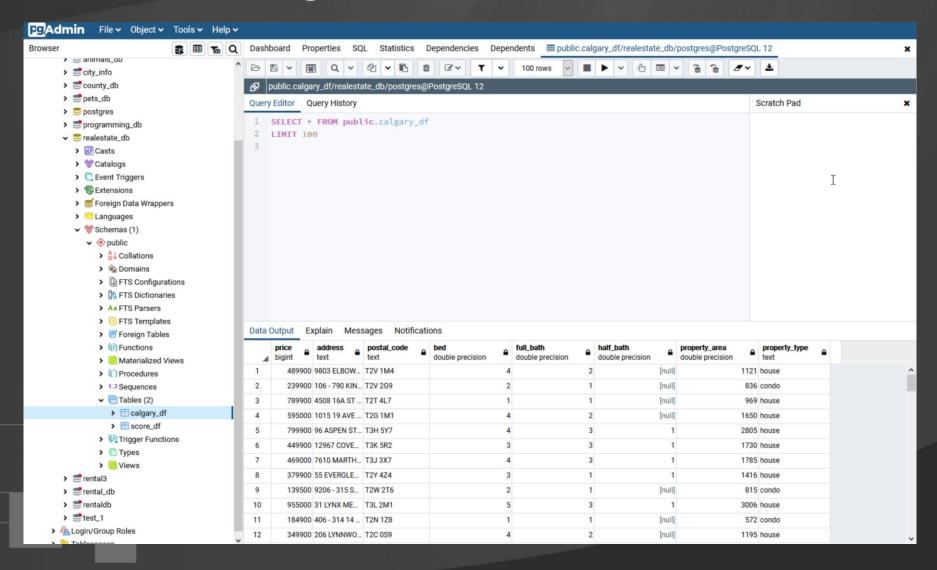
In [30]: rds_connection_string = "postgres:1@localhost:5432/realestate_db"
    engine = create_engine(f'postgresql://{rds_connection_string}')

    calgary_df.to_sql(name= "calgary_df", con=engine, if_exists="append", index=False)
    score_df.to_sql(name= "score_df", con=engine, if_exists="append", index=False)
```

PostgreSQL Database



PostgreSQL Database



Loading Data to Non-Relational Database (MongoDB)

Building connection to MongoDB and loading transformed data

MongoDB

```
In []: # Make a connection
    conn = "mongodb://localhost:27017"

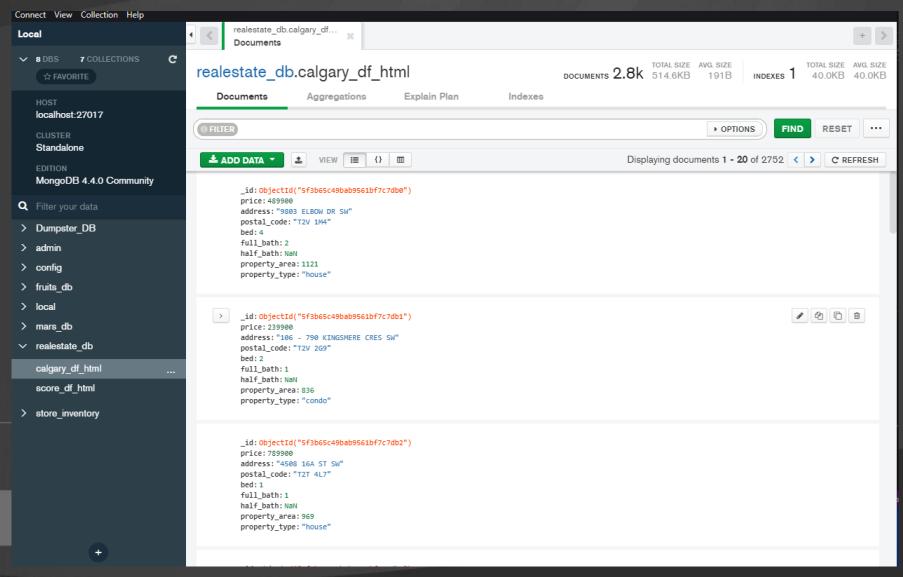
# Making a Connection with MongoClient
    client = MongoClient(conn)

# database
    db = client.realestate_db

collection = db.calgary
    calgary_dict = calgary_df.to_dict("records")
    collection.insert_many(calgary_dict)

collection = db.score
    score_dict = score_df.to_dict("records")
    collection.insert_many(score_dict)
```

MongoDB



Converted Database into a Web Based Application



Remax Calgary Realestate

Price (\$)	Address	Postal Code	Bedrooms	Bathrooms	Property Type
489900	9803 ELBOW DR SW	T2V 1M4	4 bed	2 bath	house
319900	101 - 3704 15A ST SW	T2T 4C3	2 bed	1 + 1 bath	condo
399900	25 HARVEST GLEN WAY NE	T3K 4J2	3 bed	2 bath	house
429500	32 EVERGLEN GROVE SW	T2Y 4Z3	3 bed	2 + 1 bath	house
484900	416 THORNDALE RD NW	T2K 3C5	3 bed	2 bath	house
529000	46 HARVEST GROVE CLOSE NE	T3K 4T6	4 bed	3 + 1 bath	house
379900	26 ANAHEIM PL NE	T1Y 7B3	4 bed	2 + 1 bath	house
337000	1905 - 930 6 AVE SW	T2P 1J3	1 bed	1 bath	condo
544500	52 RIVERVIEW MEWS SE	T2C 3Z8	4 bed	3 bath	house
445000	1 - 435 13 AVE NE	T2E 1C3	3 bed	2 + 1 bath	townhouse
229000	3 - 203 VILLAGE TERR SW	T3H 2L4	2 bed	2 bath	condo
409900	234 ROYAL BIRCH BAY NW	T3G 5X6	3 bed	3 bath	house
749900	111 HILLGROVE CRES SW	T2V 3K9	2 bed	1 bath	house
924900	3030 26A ST SW	T3E 2E3	3 bed	2 bath	house
574900	78 CHAPARRAL VALLEY GROVE SE	T2X 0M4	4 bed	3 + 1 bath	house
438700	203 ARBOUR STONE PL NW	T3G 5E9	3 bed	2 + 1 bath	house
369900	108 - 59 22 AVE SW	T2S 3C7	2 bed	2 bath	condo
489000	140 CITADEL CREST CIR NW	T3G 4G3	3 bed	2 + 1 bath	house
469900	25 DOUGLASBANK RISE SE	T2Z 2C5	4 bed	2 + 1 bath	house

Postal Code	Walk Score	Bike Score	Transit Score
T2V 1M4	58	61	55
T2T 4C3	53	81	42
T3K 4J2	19	59	38
T2Y 4Z3	6	30	31
T2K 3C5	61	81	53
T3K 4T6	40	60	39
T1Y 7B3	13	59	40
T2P 1J3	94	94	78
T2C 3Z8	32	70	39
T2E 1C3	82	84	54
T3H 2L4	26	67	35
T3G 5X6	47	65	36
T2V 3K9	36	53	53
T3E 2E3	56	74	48
T2X 0M4	4	50	28
T3G 5E9	15	47	49
T2S 3C7	55	91	59
T3G 4G3	23	65	38