



8/5/2018

# Data Management, Warehousing Analytics CSCI 5408

Project Report: Choosing a location for  
establishing a Business



Sogra Bilal Memon  
B00786252

Bhavya Chandrappa  
B00781097

## Table of Contents

<b>1. Summary</b>	<b>3</b>
<b>2. Business Idea, Problem statement, Value Proposition</b>	<b>3</b>
<b>3. Implementation</b>	<b>4</b>
<b>3.1. Data Source</b>	<b>4</b>
<b>3.2. Tools used</b>	<b>5</b>
<b>3.3. Planned Procedures</b>	<b>6</b>
<b>3.4. Execution of the planned procedures</b>	<b>9</b>
<b>3.5. Implementation Explained in Detail</b>	<b>10</b>
<b>3.6. Visualization</b>	<b>17</b>
<b>4. Role based work and distribution</b>	<b>27</b>
<b>5. Milestones/Sprints</b>	<b>27</b>
<b>6. Limitations of Current work</b>	<b>28</b>
<b>7. Future Work</b>	<b>28</b>
<b>8. Critical Review</b>	<b>28</b>
<b>9. References</b>	<b>29</b>

## **1. Summary**

In this project we have focused on generating visualizations, based on current demographics that could be helpful to new comers who are aspiring entrepreneurs. For simplicity we have taken a possible business use case and accordingly modeled visualizations to show which city from Nova Scotia has the highest growth potential for this venture. These visualizations of the current population along with the income, age and language distribution gives aspiring business owners an idea of his/her potential customers, their product of interest and their spending budget. We have chosen the idea of starting a senior citizen café and have visualized the current demographics to find the city where this venture has the highest growth potential.

## **2. Business Idea, Problem statement, Value Proposition**

Halifax is growing as a city and with this growth there are many people from other countries who want to migrate to Halifax. With the migration process comes a number of challenges such as visa application, renting a home, finding a school to setting up a business. While google is a great tool for all these searches it does not provide a visualization of all the data available in order for users to take immediate decisions. Users have to rely on their memory and constant google searches to make the right decisions for each of these situations. Our project aims to help newcomers set up a business with ease so that they can focus on other aspects of settling down into a new place. Our project will be working with CENSUS data to give potential business owners all the information they need about their surrounding neighborhood to come to a conclusion about which products will sell in that particular area. The CENSUS data is a record of the number of people living in each location in Canada. Demographic information such as age, gender, marital status and languages spoken are also available for people living in each location. This makes the CENSUS data the perfect dataset for potential business owners to look through before they take a plunge at their new business idea.

We aim to make the CENSUS data available to anyone who is looking to validate the feasibility of setting up their venture at a specific location. Our project will also help users find the best possible location for his or her venture based on the requirements for business success provided by the user. The user will be provided with data visualizations that will provide clear indications on which locations are better for the business as compared to others based on the parameters selected by the user. These visualizations can also be used by the user to develop business proposals and strategies that are not based on intuition but on actual data that can be presented to investors as evidence that the product or service will work.

Starting a business in Halifax is a goal of many Halifax citizens. Many times the users get stuck in locating a place to set up their business. There is no uncomplicated way to make the users' search successful. Our project aides the users to locate the suitable location by using data analytics concepts and tools. The data is and has been available for a long time. This project

aims to make it available to public in a way that it is easy to use and meets the demands of providing some certainty to a startup that is in an uncertain risky environment. Many startups fail every year not necessarily because the venture idea was bad but because it was executed at the wrong location. The location and timing are key to any business. With this solution, users will be able to make better choices based on data and quantitative comparison between locations.[3]

### 3. Implementation

The implementation of this project has been achieved through a series of sprint. The sprints have been discussing in detail in section 5. In this section we discuss the data set used and the procedures followed to clean, process, visualize and analyze the data.

#### 3.1. Data Source

The Dataset being used in this project has been taken from CENSUS Canada. This dataset consists of all the information gathered by CENSUS in 2016. CENSUS is a survey that takes place every five years. The dataset taken contains information on population, income and language distribution for the entire Canada. These records have been taken for each of the counties in each of the provinces which caused it to occupy 84mb of memory. The data was highly unstructured and required the extraction and transformation of relevant features. The data occupies 84 Mb of memory space and has a basic schema as shown below.

P	<b>Id</b>	<b>City</b>	<b>Feature</b>	<b>Feature value</b>
	1	Halifax	Total population	69084364
	1	Halifax	Population over age of 65	14
	1	Halifax	Population age range 15 - 64	143536
	1	Halifax	Population percentage age 65 and over	40
	1	Halifax	Population percentage age range 15 to 64	256678
	1	Halifax	Total population with Income	23546576
	1	Halifax	Population in income range under 10,000	23465787
	1	Halifax	Population that speaks Hindi	23456

Figure 1: Work Structure Breakdown

### 3.2. Tools used

In this project we have used a variety of tools to store, clean, process, transform and perform visualization on the dataset chosen. Some of the tools and languages used include: Python, AWS EC2 instance, Spark, MySQL and Tablue. The details of how each of the tools that were used has been described in section 3.3.

### 3.3. Planned Procedures

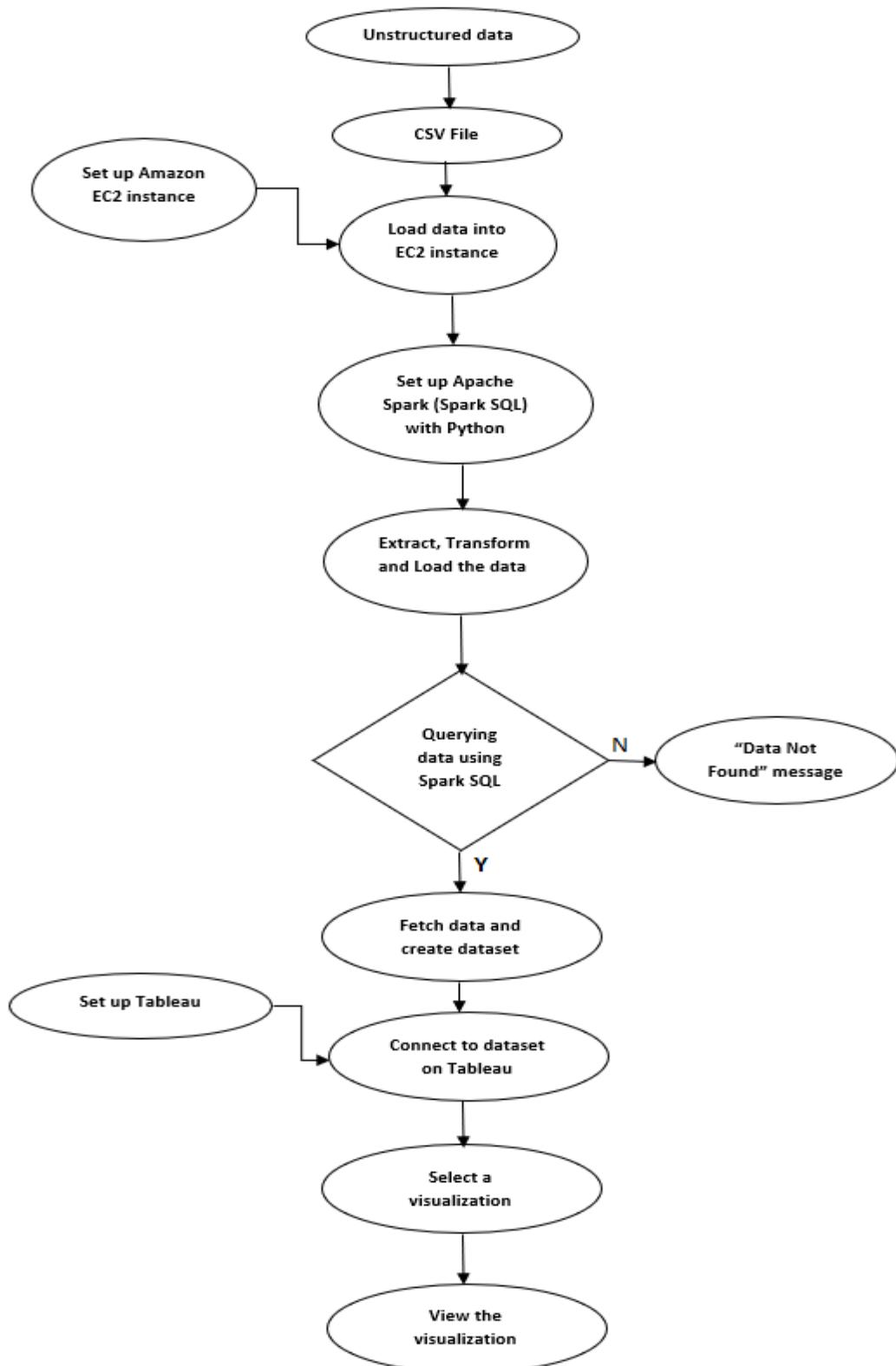


Figure 1: Work Structure Breakdown

In this project we planned to collect, clean and manipulate the CENSUS the following way.

### Data Extraction

The CENSUS data seemed to be well structured in the form of tables, rows and columns. However, it is not available in the form of a csv file.

The diagram illustrates the breakdown of a large census dataset into smaller, more manageable tables. The original dataset is shown as a large table on the left, which is then broken down into three smaller tables on the right. Arrows point from specific rows in the original table to their corresponding derived tables.

Id	City	Feature	Feature value
1	Halifax	Total population	69084364
1	Halifax	Population over age of 65	14
1	Halifax	Population age range 15 - 64	143536
1	Halifax	Population percentage age 65 and over	40
1	Halifax	Population percentage age range 15 to 64	256678
1	Halifax	Total population with Income	23546576
1	Halifax	Population in income range under 10,000	23465787
1	Halifax	Population that speaks Hindi	23456

Id	City	Feature	Feature value
1	Halifax	Total population	69084364
1	Halifax	Population over age of 65	14
1	Halifax	Population age range 15 - 64	143536
1	Halifax	Population percentage age 65 and over	40
1	Halifax	Population percentage age range 15 to 64	256678

Id	City	Feature	Feature value
1	Halifax	Total population with Income	23546576
1	Halifax	Population in income range under 10,000	23465787

Id	City	Feature	Feature value
1	Halifax	Population that speaks Hindi	23456

Figure 1: Work Structure Breakdown

## Data Cleaning

Once the CSV file is ready the data can be cleaned to remove any null values as well as any unnecessary columns. Depending on the quality of the data it may also have to be normalized.

## Setting up an EC2 instance and loading the data

This project will run on an EC2 instance which will be setup and configured with an Ubuntu virtual machine and the required software packages such as Python and apache spark will be installed. The CSV files prepared can then be loaded into the EC2 instance.

## After ETL processes

Further ETL processes will be carried out on the csv file so that the data can be loaded in apache spark. This will further remove any noise in the data.

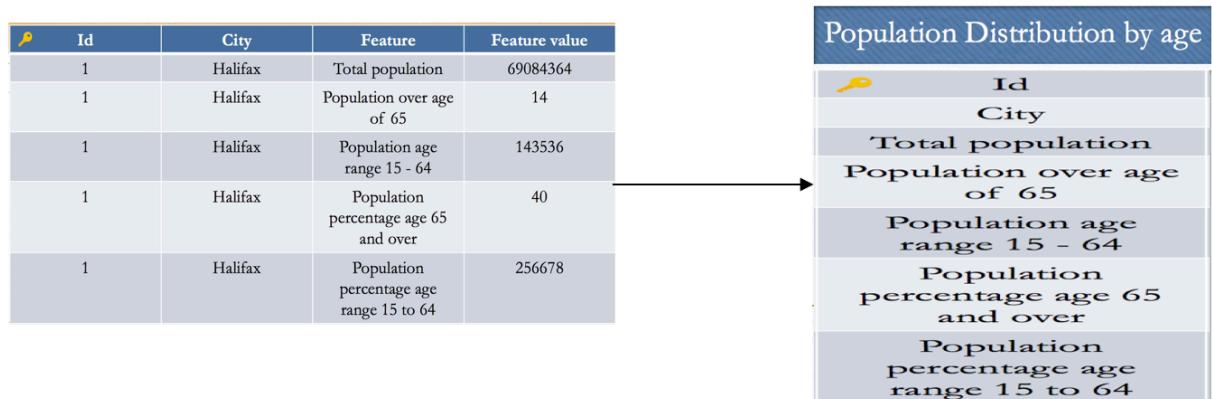


Figure 1: Work Structure Breakdown

## Querying the data

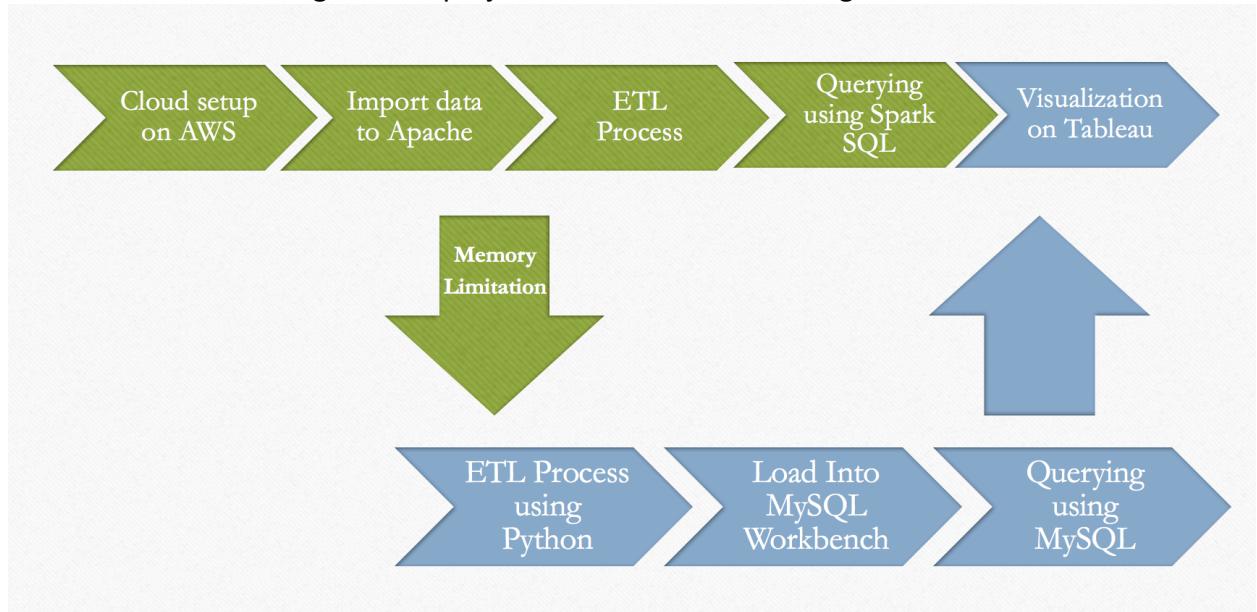
Queries will be generated based on user's requirements. These queries will be executed by Spark SQL to retrieve the required data. If there is no data available that matches the search query an error message will be displayed. On the other hand, if the query is executed successfully the data is fetched and saved.

Connect to tableau

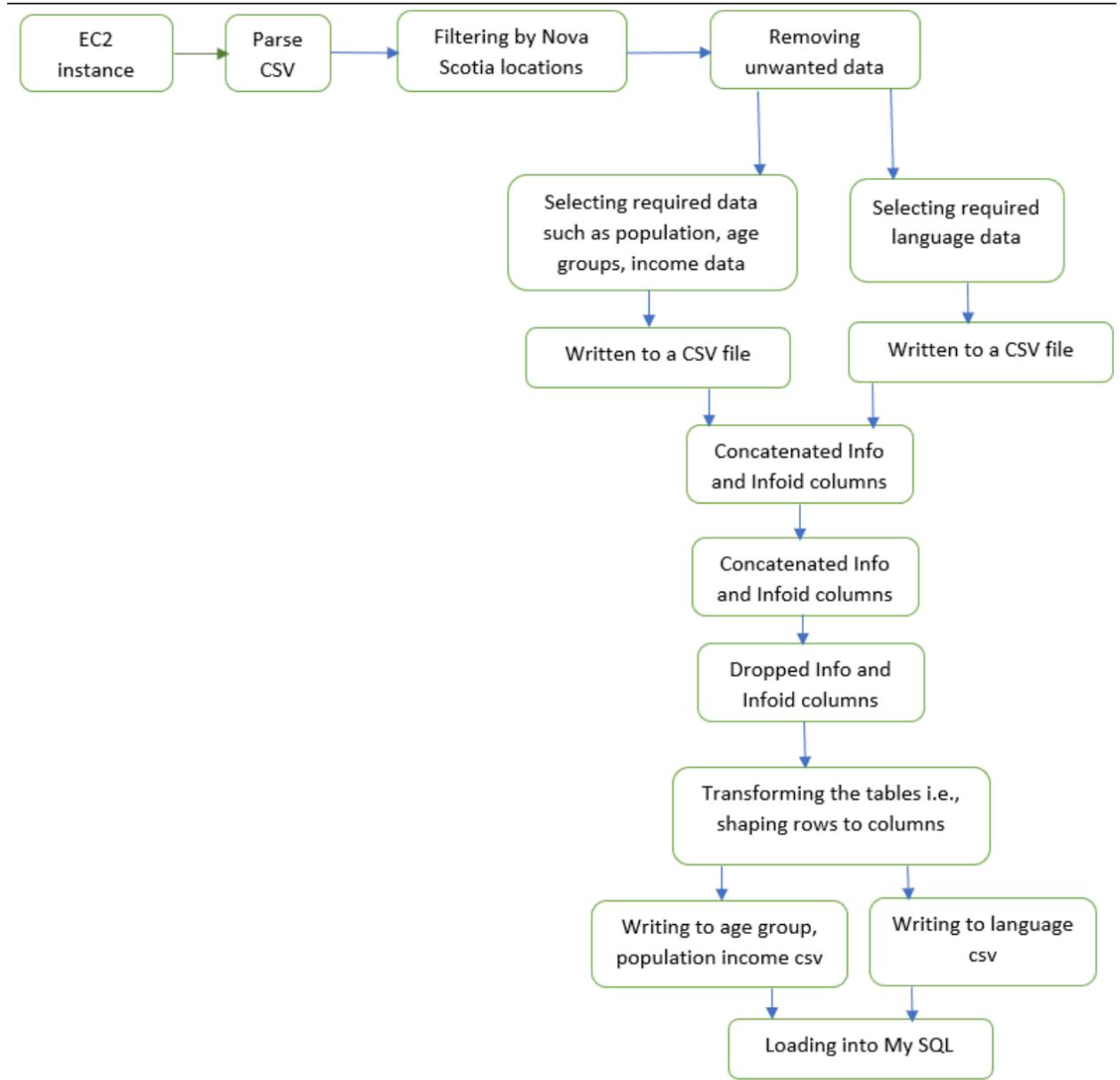
The dataset resulting from the query is connected to tableau and the user is given options to choose the type of visualization he/she would like. Once the type of graph is selected tableau generates and presents the visualization generated.

### 3.4. Execution of the planned procedures

This project required us to collect, clean and manipulate the CENSUS data to be able to achieve the desired insights. This project will follow the following breakdown structure.



### 3.5. Implementation Explained in Detail



#### Data Extraction

The CENSUS data turned out to be more unstructured than we expected. The schema made it very hard to query especially join functions. We had to extract the required data into separate data frames. Where each data frame will have a similar set of data. In the example below we show the complete dataset being split into the population, income and language data frames. We also realized that the data set was too large and decided to only take the records that belonged to Nova Scotia instead of for the entire Canada.

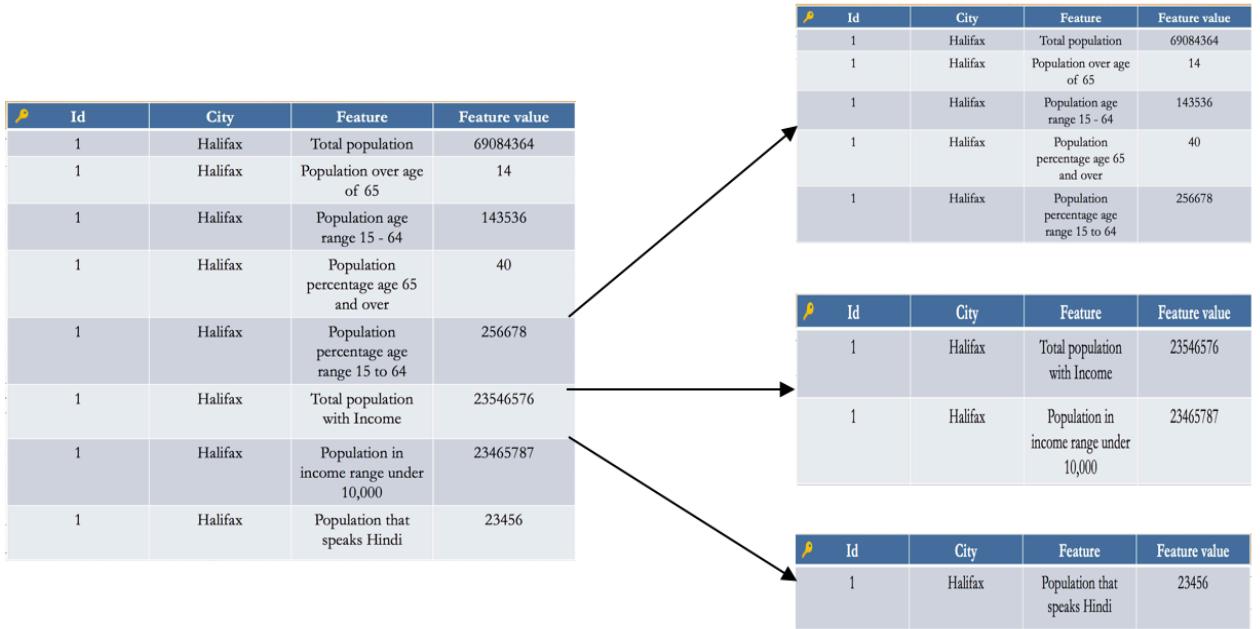


Figure 1: Work Structure Breakdown

### A. Data Extraction with Spark

As mentioned in the proposal we planned to use spark for the ETL processes and Spark SQL for querying. However, we ran into memory errors while trying to process our entire dataset. As a result a python scripts where run on the deployed EC2 instance to clean and extract the required data.

```
PuTTY (inactive)
Welcome to
   _/\_  _/\_  _/\_  _/\_  _/\_  _/\_
  / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \
 version 2.3.0

Using Python version 3.4.3 (default, Nov 28 2017 16:41:13)
SparkSession available as 'spark'.
>>> import csv
>>> from pyspark import SparkContext, SparkConf
>>> from pyspark.sql import SQLContext
>>> sc.stop()
>>> sc = SparkContext("spark://ip-172-31-34-76.us-east-2.compute.internal:7077","Simple App")
>>> sqlContext = SQLContext(sc)
>>> df = sqlContext.read.csv("census.csv")
2018-07-27 14:40:09 WARN ObjectStore:568 - Failed to get database global_temp, returning NoSuchObject
Exception: [Stage 0]
[Stage 0]>                                         (0 + 1) / 1]quit
[Stage 0]>                                         (0 + 1) / 1]2018-07-27 14:44:42 ER
ROR TaskSchedulerImpl:70 - Lost executor 0 on 172.31.34.76: Remote RPC client disassociated. Likely due to containers exceeding thresholds, or network issues. Check driver logs for WARN messages.
2018-07-27 14:44:43 WARN TaskSetManager:66 - Lost task 0.0 in stage 0.0 (TID 0, 172.31.34.76, executor 0): ExecutorLostFailure (executor 0 exited caused by one of the running tasks) Reason: Remote RPC client disassociated. Likely due to containers exceeding thresholds, or network issues. Check driver logs for WARN messages.
[Stage 0]>                                         (0 + 0) / 1]quit
exit
<CTraceback (most recent call last):
  File "<stdin>", line 1, in <module>
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/pyspark/sql/readwriter.py", line 439, in
csv
    return self._df(self._jreader.csv(self._spark, sc._jvm.PythonUtils.toSeq(path)))
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/lib/py4j-0.10.6-src.zip/py4j/java_gateway
.py", line 1158, in _call
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/lib/py4j-0.10.6-src.zip/py4j/java_gateway
.py", line 908, in send_command
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/lib/py4j-0.10.6-src.zip/py4j/java_gateway
.py", line 1055, in send_command
  File "/usr/lib/python3.4/socket.py", line 374, in readline
    return self._sock.recv_into(b)
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/pyspark/context.py", line 234, in signal_
handler
    raise KeyboardInterrupt()
KeyboardInterrupt
>>> df = sqlContext.read.format("csv").option("header","true").load("census.csv")
Java HotSpot(TM) 64-Bit Server VM warning: INFO: os::commit_memory(0x00007f9c76742000, 262144, 0) fail
ed; error='Cannot allocate memory' (errno=12)
#
# There is insufficient memory for the Java Runtime Environment to continue.
#
Windows PowerShell Type here to search 91% ENG 1058 AM 2018-07-29
```

```
PuTTY (inactive)
py4j.protocol.Py4JNetworkError: Answer from Java side is empty

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/lib/py4j-0.10.6-src.zip/py4j/java_gateway
.py", line 908, in send_command
    response = connection.send_command(command)
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/lib/py4j-0.10.6-src.zip/py4j/java_gateway
.py", line 990, in start
    self.socket.connect((self.address, self.port))
ConnectionRefusedError: [Errno 111] Connection refused
ERROR:py4j:java gateway:An error occurred while trying to connect to the Java server (127.0.0.1:58467)
Traceback (most recent call last):
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/lib/py4j-0.10.6-src.zip/py4j/java_gateway
.py", line 852, in _get_connection
    connection = self._deque.pop()
IndexError: pop from an empty deque

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/lib/py4j-0.10.6-src.zip/py4j/java_gateway
.py", line 990, in start
    self.socket.connect((self.address, self.port))
ConnectionRefusedError: [Errno 111] Connection refused
ERROR:py4j:java gateway:An error occurred while trying to connect to the Java server (127.0.0.1:58467)
Traceback (most recent call last):
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/lib/py4j-0.10.6-src.zip/py4j/java_gateway
.py", line 852, in _get_connection
    connection = self._deque.pop()
IndexError: pop from an empty deque

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/lib/py4j-0.10.6-src.zip/py4j/java_gateway
.py", line 990, in start
    self.socket.connect((self.address, self.port))
ConnectionRefusedError: [Errno 111] Connection refused
ERROR:py4j:java gateway:An error occurred while trying to connect to the Java server (127.0.0.1:58467)
Traceback (most recent call last):
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/lib/py4j-0.10.6-src.zip/py4j/java_gateway
.py", line 852, in _get_connection
    connection = self._deque.pop()
IndexError: pop from an empty deque

During handling of the above exception, another exception occurred:

Traceback (most recent call last):
Windows PowerShell Type here to search 91% ENG 1059 AM 2018-07-29
```

```

PuTTY (inactive)
Using Python version 3.4.3 (default, Nov 28 2017 16:41:13)
SparkSession available as 'spark'
>>> from pyspark import SparkContext, SparkConf
>>> from pyspark.sql import SQLContext
>>> from pyspark.sql.types import *
>>> import csv
>>> sc = SparkContext("spark://ip-172-31-34-76.us-east-2.compute.internal:7077","Simple App")
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/pyspark/context.py", line 115, in __init__
    SparkContext._ensure_initialized(self, gateway=gateway, conf=conf)
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/pyspark/context.py", line 296, in ensure_initialized
    callsite.function, callsite.file, callsite.linenum)
ValueError: Cannot run multiple SparkContexts at once: existing SparkContext(app=PySparkShell, master=spark://ip-172-31-34-76.us-east-2.compute.internal:7077) created by <module> at /home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/pyspark/shell.py:45
>>> sc = SparkContext("spark://ip-172-31-34-76.us-east-2.compute.internal:7077","Simple App")
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/pyspark/context.py", line 115, in __init__
    SparkContext._ensure_initialized(self, gateway=gateway, conf=conf)
  File "/home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/pyspark/context.py", line 296, in ensure_initialized
    callsite.function, callsite.file, callsite.linenum)
ValueError: Cannot run multiple SparkContexts at once: existing SparkContext(app=PySparkShell, master=spark://ip-172-31-34-76.us-east-2.compute.internal:7077) created by <module> at /home/ubuntu/server/spark-2.3.0-bin-hadoop2.7/python/pyspark/shell.py:45
>>> sc.stop()
>>> sc = SparkContext("spark://ip-172-31-34-76.us-east-2.compute.internal:7077","Simple App")
>>> sqlContext = SQLContext(sc)
>>> df = sqlContext.read.csv("census.csv")
2018-07-27 06:22:24 WARN ObjectStore:6666 - Version information not found in metastore. hive.metastore.schema.verification is not enabled so recording the schema version 1.2.0
2018-07-27 06:22:29 WARN ObjectStore:568 - Failed to get database default, returning NoSuchObjectException
2018-07-27 06:22:50 WARN ObjectStore:568 - Failed to get database global_temp, returning NoSuchObjectException
exit()
quit()
[Stage 0:>
exit():

```

(0 + 1) / 1]

## B. Data Extraction with Python

When Spark started to give memory errors we switched to using python scripts to extract the features we needed and remove any columns and rows with null values. The features we choose to extract were based on the visualizations we wanted to perform to determine the best city to launch our business idea “The Senior Café”. We decided to extract features such as the percentage of the senior citizen population in each city.

```

ubuntu@ip-172-31-20-159:~ File: pytnnew.py
[1]  GNU nano 2.2.6
import csv
import pandas as pd
import io
import numpy as np

#@infoline = csv.reader(open('census.csv','r'))
cols=['Year','Geocode','Geolevel','Geoname','Gnr','Gnrlf',
'Dataflag','Altgeo','Info','Infoid','Notes','Total','Male','Female']
dt cols=[Year:int,Geocode:int,Geolevel:float,Geoname:str,Gnr:float,Gnrlf:float,
Dataflag:float,Altgeo:float,Info:str,Infoid:int,Notes:float,Total:float, Male:float,Female:float]
##df=pd.DataFrame(np.empty(0,dtype=dt_cols))
try:
    df = pd.read_csv('newfile.csv', skiprows=1, index_col=False, names=cols, low_memory=False)## dtype=dt_cols)
except CParserError:
    print("Something wrong the file")
##df = pd.read_csv('newfile.csv', skiprows=1, index_col=False, names=cols, low_memory=False)## dtype=dt_cols)
##time.sleep(3)
print("Imported data-----")
print(df.head(3))
print("Number of rows before filter: ",df.shape[0])
print("Filtering nova scotia cities/towns-----")
filter_list = ['Amherst','Antigonish','Bedford','Bridgewater','Digby','Glace Bay','Inverness','New Waterford','Halifax','Kentville','Middleton','New Glasgow','Sydney Mines','Shelburne','Cape Breton - Sydney','Tiruro','Wolfville','Yarmouth','Chester','Berwick','Hantsport','Lunenburg','Springhill','Lake Echo','Hayes Subdivision','Brookside','Still Water Lake','Enfield - Lantz','English Corner','Port Williams','Howie Centre','Pictou']
##df[df['Geoname'].isin(filter_list)]
df1=df.loc[df['Geoname'].isin(filter_list)]
print("Number of rows after filter: ",df1.shape[0])
##df.columns = ['code', 'geo_name', 'lang', 'lang_id', 'total', 'men', 'women']
##df.apply(lambda x: pd.lib.infer_dtypes(x.values))

#for i, row in df.iterrows():
#    # print(row)
#    # df=df['Info'].replace(['-',' ','_',' '],regex = True)
#    # df=df['Geoname'].replace(['Amherst'],['Amherst'],regex = True)

#df2=pd.concat([d,df],axis=1)
print("Dropped data-----")
df2=df1.drop(['Year','Geolevel','Gnr','Gnrlf','Dataflag','Altgeo','Notes','Male','Female'], axis=1)
df2.to_csv('dropdata.csv',index=False,encoding='utf-8')
print(df2.head(3))
print("Infoid 34 to 38 data(distribution of pop)----")
df3 = df2[df1.Infoid > 36]

```

Get Help WriteOut Read File Prev Page Cut Text Cur Pos  
Exit Justify Where Is Next Page Uncut Text To Spell  
Type here to search 01:07 ENG 10:38 PM 2018-08-05

```

ubuntu@ip-172-31-20-159:~ File: pytnnew.py
[1]  GNU nano 2.2.6
print("kail data-----")
print(df14.tail(3)) ##df14
print("Infoid 121-375 data(languages)----")
df15=df14[df14.Infoid > 120]
print("Head data-----")
print(df15.head(3))
print("Tail data-----")
print(df15.tail(3))
df16=df15[df15.Infoid <376]
print("Head data-----")
print(df16.head(3))
print("Tail data-----")
print(df16.tail(3))
df17=df16.loc[df16['Info'].str.contains('languages', regex=False, case=False, na=False)] ##df17
for i, row in df3.iterrows():
    df4=df3['Info'].replace(['%'], ['percentage'], regex = True)
##df7=pd.concat([df4,df6],axis=0)
df18=df4.append(df6)
df19=df18.append(df8)
df20=df19.append(df10)
df20.to_csv('withoutlang.csv',index=False,encoding='utf-8')
df21=df14.append(df12)
df22=df21.append(df17)
df22.to_csv('lang.csv',index=False,encoding='utf-8')
print("Writing to csv----")
##print(df23.head(3))
##df23.to_csv('conatenewnew.csv',index=False,encoding='utf-8')
##df2.replace('...',' ', regex=True)
##df2.fillna(0)
print("Reshaping-----")
##df16['Info_Infoid']=df16['Info']+ ' '+df16['Infoid'].astype(str)
df20.loc[:, 'Info_Infoid']=df20.Infoid.map(str) + " " + df20.Info
df22.loc[:, 'Info_Infoid']=df22.Infoid.map(str) + " " + df22.Info
##df16.to_csv('conact.csv',index=False,encoding='utf-8')
df23=df20.drop(['Infoid'],axis=1)
df26=df22.drop(['Infoid'],axis=1)
df24=df23.drop(['Info'],axis=1)
df27=df26.drop(['Info'],axis=1)
df28=df27.pivot_table(['Total'], ['Geocode','Geoname'], ['Info_Infoid'], aggfunc='first')
df29=df28.pivot_table(['Total'], ['Geocode','Geoname'], ['Info_Infoid'], aggfunc='first')
##df18=df18.drop(['Infoid'],axis=1)
print("Writing to pivoted csv----")
print(df25.head(3))
df25.to_csv('pivot.csv',index=True,encoding='utf-8')

```

Get Help WriteOut Read File Prev Page Cut Text Cur Pos  
Exit Justify Where Is Next Page Uncut Text To Spell  
Type here to search 01:07 ENG 10:39 PM 2018-08-05

## Data Cleaning

Once the CSV files were ready the data was cleaned to remove any null values as well as any unnecessary columns. This dataset had a lot of null values that were removed.

## Setting up an EC2 instance and loading the data

The EC2 instance was setup opting for the highest memory possible with the free student account. The required programming languages and libraries were installed. Some of the packages installed include Python, pandas library, pip library and MySQL. The dataset was then loaded into the EC2 instance. Python scripts were then written and applied to the loaded dataset to generate extracted csv files.

## ETL processes

After Extracting the required data each of the CSV files were transformed using python scripts to change the features recorded in rows to columns. The change in schema is shown below.

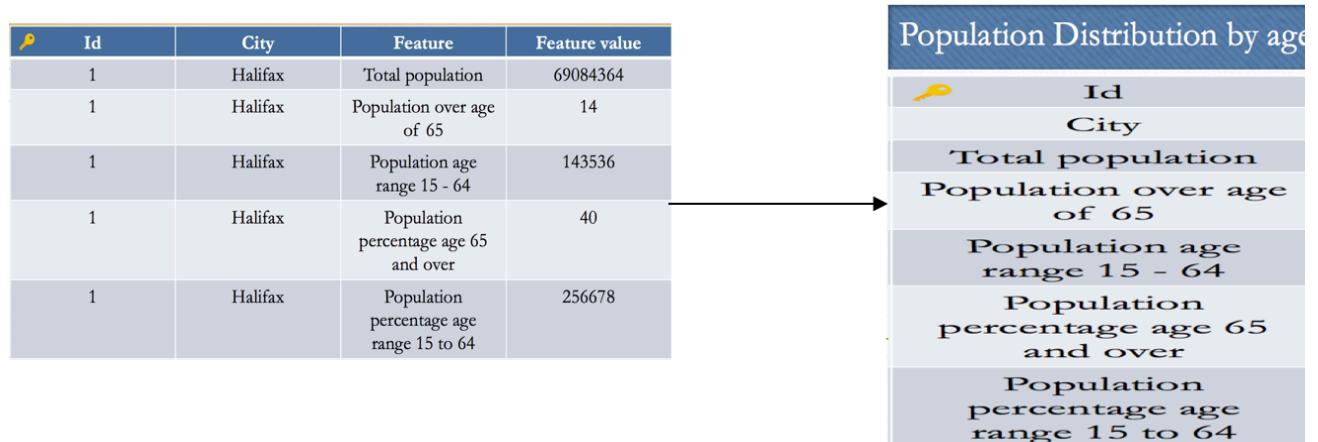


Figure 1: Work Structure Breakdown

## Querying the data

The dataset has been transformed from an unstructured dataset to a structured schema. This schema was then loaded into MySQL workbench to allow the data to be queried. The queries were chosen based on our business idea. For the future of this application the queries would be generated dynamically to generate current market visualization for any type of business desired. The query takes a CSV as an input and gives a much smaller more clean csv as a result that is then exported to tableau for visualization. The database schema after ETL processes and loading into MySQL workbench is shown below.

**pivot**

- Geocode INT(11)
- Geoname TEXT
- a10\_0\_to\_4\_years INT(11)
- b11\_5\_to\_9\_years INT(11)
- c12\_10\_to\_14\_years INT(11)
- d13\_15\_to\_64\_years INT(11)
- e14\_15\_to\_19\_years INT(11)
- f15\_20\_to\_24\_years INT(11)
- g16\_25\_to\_29\_years INT(11)
- h17\_30\_to\_34\_years INT(11)
- i18\_35\_to\_39\_years INT(11)
- j19\_40\_to\_44\_years INT(11)
- k20\_45\_to\_49\_years INT(11)
- l21\_50\_to\_54\_years INT(11)
- m22\_55\_to\_59\_years INT(11)
- n23\_60\_to\_64\_years INT(11)
- o24\_65\_years\_and\_over INT(11)
- p37\_65\_years\_and\_over DOUBLE
- q38\_85\_years\_and\_over DOUBLE
- r711\_Percentage\_with\_after-tax\_income DOUBLE
- s712\_Under\_10000\_(including\_loss) INT(11)
- t713\_10000\_to\_19999 INT(11)
- u714\_20000\_to\_29999 INT(11)
- v715\_30000\_to\_39999 INT(11)
- w716\_40000\_to\_49999 INT(11)
- x717\_50000\_to\_59999 INT(11)
- y718\_60000\_to\_69999 INT(11)
- z719\_70000\_to\_79999 INT(11)
- aa720\_80000\_and\_over INT(11)
- bb721\_80000\_to\_89999 INT(11)

10 more...

**Indexes**

**langpivot**

- Geocode INT(11)
- Geoname TEXT
- 112\_Total\_Mother\_tongue INT(11)
- 113\_Single\_responses INT(11)
- 114\_Official\_languages INT(11)
- 115\_English INT(11)
- 116\_French INT(11)
- 117\_Nonofficial\_languages INT(11)
- 118\_Aboriginal\_languages INT(11)
- 119\_Algonquian\_languages INT(11)
- 121\_CreeMontagnais\_languages INT(11)
- 132\_Eastern\_Algonquian\_languages INT(11)
- 135\_OjibwayPotawatomi\_languages INT(11)
- 140\_Algonquian\_languages\_nie INT(11)
- 141\_Athabaskan\_languages INT(11)
- 142\_Northern\_Athabaskan\_languages INT(11)
- 152\_SlaveyHare\_languages INT(11)
- 156\_Tahltan\_languages INT(11)
- 159\_Tutchone\_languages INT(11)
- 162\_Athabaskan\_languages\_nie INT(11)
- 164\_Inuit\_languages INT(11)
- 167\_Inuit\_languages\_nie INT(11)
- 168\_Iroquoian\_languages INT(11)
- 172\_Iroquoian\_languages\_nie INT(11)
- 175\_Salish\_languages INT(11)
- 184\_Salish\_languages\_nie INT(11)
- 185\_Siouan\_languages INT(11)
- 188\_Siouan\_languages\_nie INT(11)
- 190\_Tsimshian\_languages INT(11)
- 194\_Wakashan\_languages INT(11)

61 more...

**Indexes**

Connect to tableau

The dataset resulting from the query is connected to tableau and visualizations are generated by specifying the parameters used for the visualization.

## 3.6. Visualization

Query 1: List the population percentage of senior citizens for each city in Nova Scotia

The screenshot shows the MySQL Workbench interface with the following details:

- File Bar:** MySQL Workbench, bubb, File, Edit, View, Query, Database, Server, Tools, Scripting, Help.
- Navigator:** MANAGEMENT (Server Status, Client Connections, Users and Privileges, Status and System Variables, Data Export, Data Import/Restore), INSTANCE (Startup / Shutdown, Server Logs, Options File), PERFORMANCE (Dashboard, Performance Reports, Performance Schema Setup).
- Schemas:** Filter objects, Tables (langpivot, pivot), Views, Stored Procedures, Functions.
- Information:** Table: langpivot, Columns: Geocode, Geoname.
- Query Editor:** SQL File 8\* contains the query:

```
select Geoname, p37_65_years_and_over from dwhproj.pivot
group by Geoname, p37_65_years_and_over;
```
- Result Grid:** Shows the results of the query:

Geoname	p37_65_years_and_over
Amherst	23.9
Antigonish	25.6
Bedford	27.6
Berwick	34.2
Bridgewater	27
Brookside	12.5
Cape Breton - Svdnev	23.2
Chester	39.4
Diaiby	31.1
Enfield - Lantz	12.9
Endish Corner	12.6
Glace Bay	23.7
- Action Output:** Shows the log entry: # 46 22:01:22 DEALLOCATE PREPARE stmt.

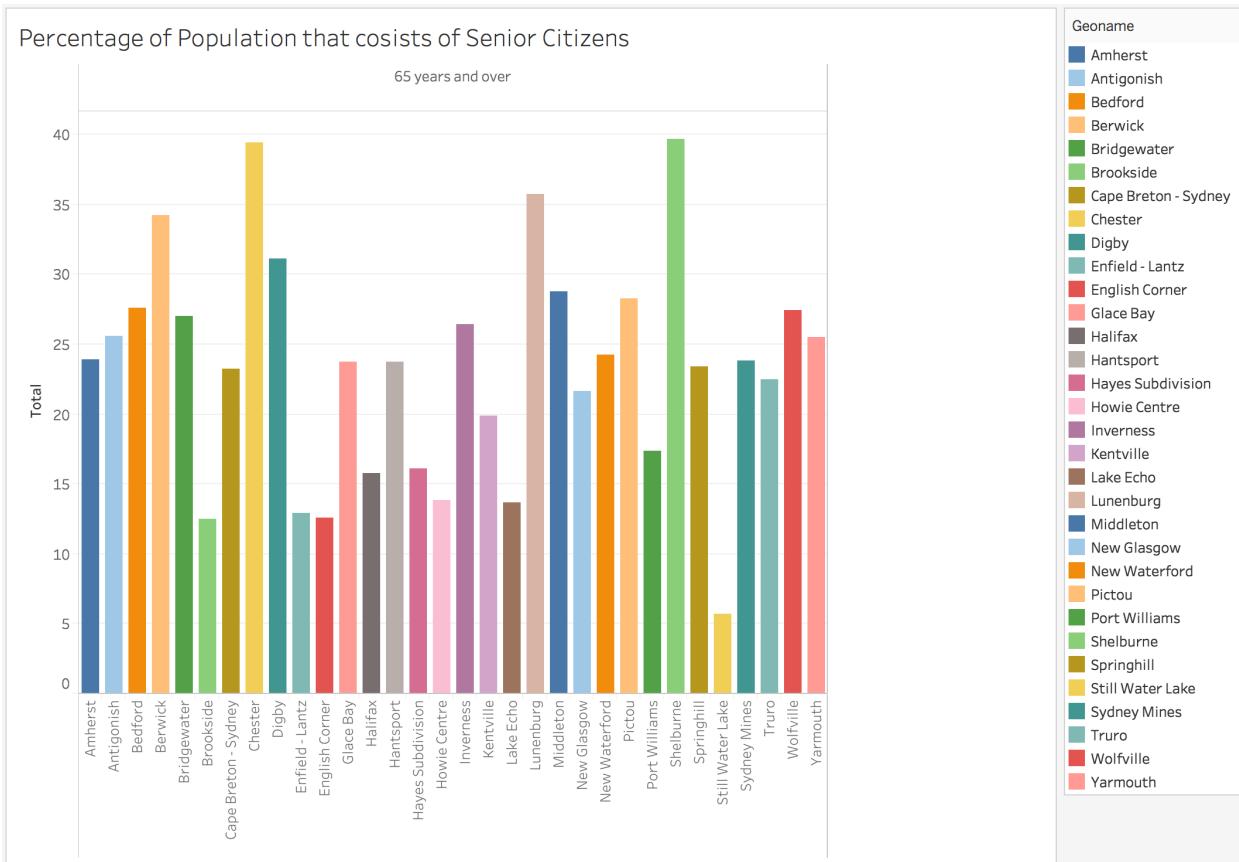


Figure 1: Work Structure Breakdown

According to this visualization Shelburne has the highest percentage of senior citizens. This information would be helpful to someone looking to start a business as they could decide to start in Shelburne based on the given information about the current market scenario. This visualization could be used in a business pitch to support the business owner's decision while presenting to a potential investor.

However, this visualization below shows that even though Shelburne has a higher percentage of the elderly population, The total population of Halifax is higher and hence Halifax has a larger population of senior citizens as compared to Shelburne.

Population Of A City And Its Age Distribution													Population of city
Info	InfoId	Halifax	Geoname										6 316,690
			Cape Breton - Sydney	Truro	Glace Bay	New Glasgow	Sydney Mines	Kentville	Bridgewater	Amherst	Yarmouth	New Waterford	
15 to 64 years	13	220,945	18,785	14,435	10,985	11,700	8,000	7,755	5,115	5,830	4,410	4,630	
	36		70	63	63	63	62	64	60	61	61	63	
65 years and over	24	11.	50,055	6,950	5,170	4,160	4,025	3,045	2,410	2,300	2,285	1,840	1,780
			16	23	23	24	22	24	20	27	24	26	24
Total - Age groups and average age of the population - 100% data	8	316,690	29,910	22,950	17,545	18,655	12,820	12,085	8,530	9,550	7,215	7,345	

Figure 1: Work Structure Breakdown

## Query 2:

The low income of the seniors helps us to decide on where to start a business for the seniors since they should not feel the business is expensive from their perspective. Below, the graph depicts the low income of the seniors versus the location. We can say that Halifax has the highest number of low income seniors who are 65 years and above.

MySQL Workbench

bubb ×

File Edit View Query Database Server Tools Scripting Help

Navigator MANAGEMENT INSTANCE SCHEMAS INFORMATION Table: pivot

Query 1 SQL File 1\* SQL File 2\* pivot - Table SQL File 5\* SQL File 6\*

```
1 select Geoname,ii851_65_years_and_over
2 from dwhproj.pivot
3 group by Geoname,ii851_65_years_and_over DESC;
```

Result Grid Filter Rows: Export: Wrap Cell Contents:

Geoname	ii851_65_years_and_over
Amherst	2055
Antigonish	1055
Bedford	600
Berwick	595
Bridgewater	2120
Brookside	175
Cape Breton - Sydney	6225
Chester	490
Digby	560
Enfield - Lantz	875
English Corner	90
Glace Bay	3950
Halifax	46760
Hantsport	360
Hayes Subdivision	180
Howie Centre	240
Inverness	260
Kentville	23900
Lake Echo	345
Lunenburg	590
Middleton	330
New Glasgow	3760
New Waterford	1700
Pictou	630

Action Output

Message OK

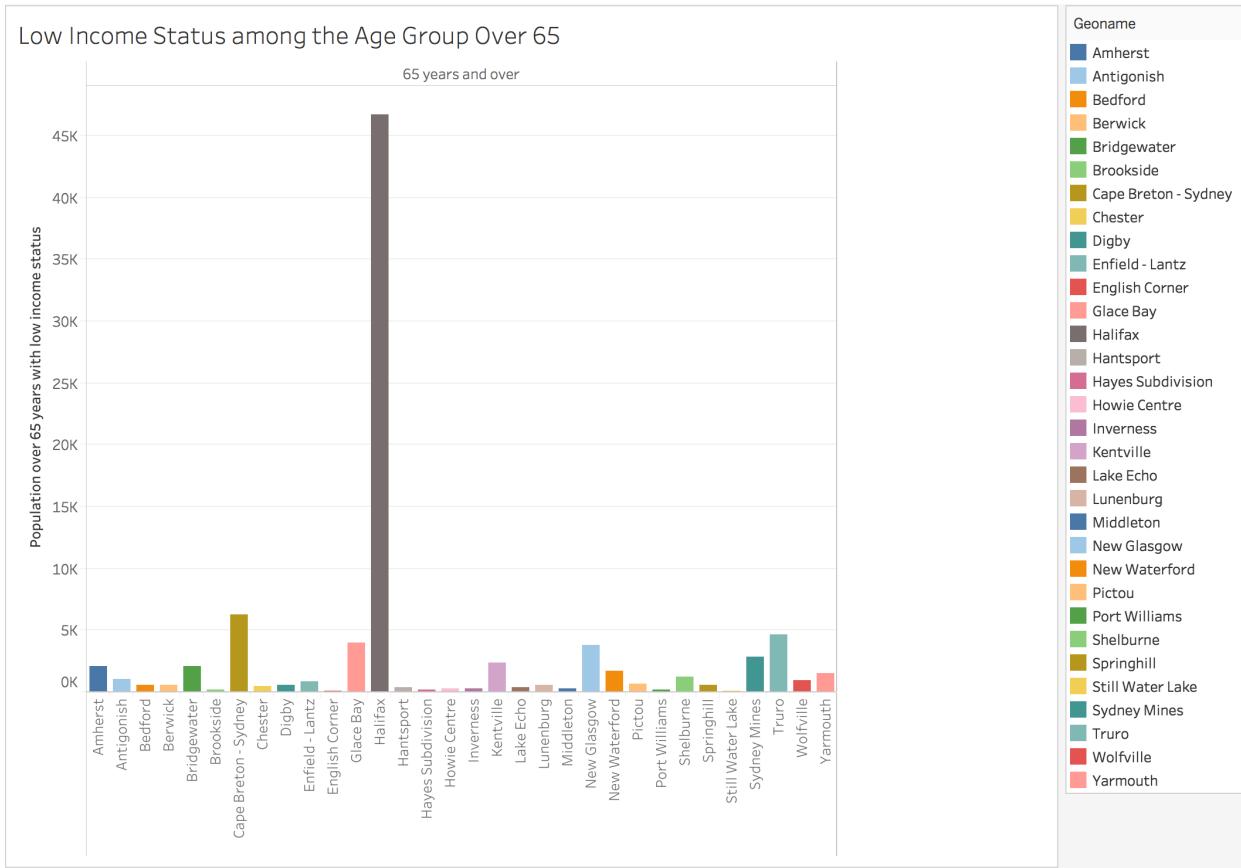


Figure 1: Work Structure Breakdown

## Languages:

Languages play an important role in our lives. Though English and French are the most spoken by the population all over Canada. There are people who still converse in mother tongue. Conversing in mother tongue makes one feel comfortable. Starting a business for seniors and letting them talk in the language they want to improves the business. They might find more people who speaks their language on daily basis. The business may also run on a specific theme where a specific ethnicity has been followed by more people in a location compared to other locations.

Below are the graphs that conveys most number of languages being spoken by Halifax population. The first graph shows the non-aboriginal languages being spoken in the Halifax where the value is represented by its size. Indo-European is spoken by most and Afro-asiatic languages come in second. If the business is targeted towards these specific languages, the number of customers who speak the same language may like the idea of meeting more people in the same location. The second graph depicts the aboriginal languages such as Algonquian languages, Easter Algonquian languages and Mi'kmaq languages

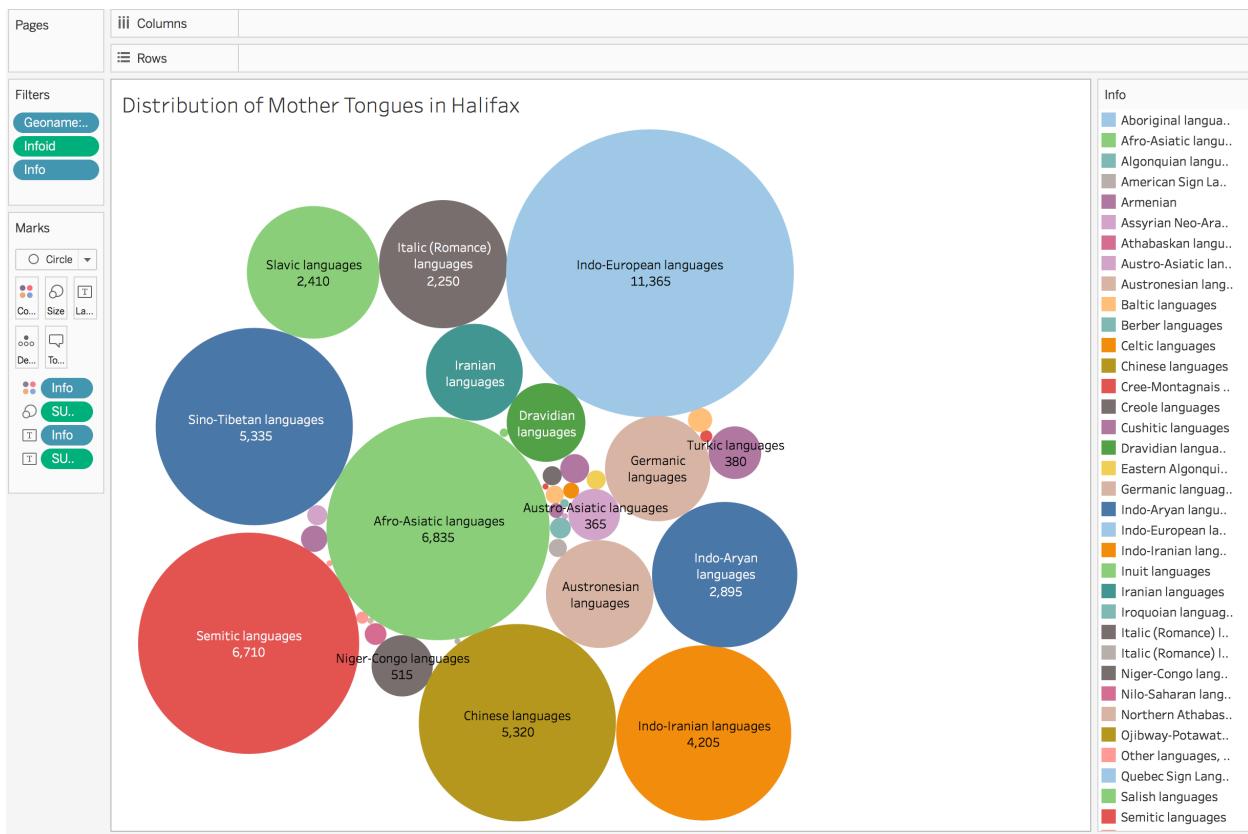


Figure 1: Work Structure Breakdown

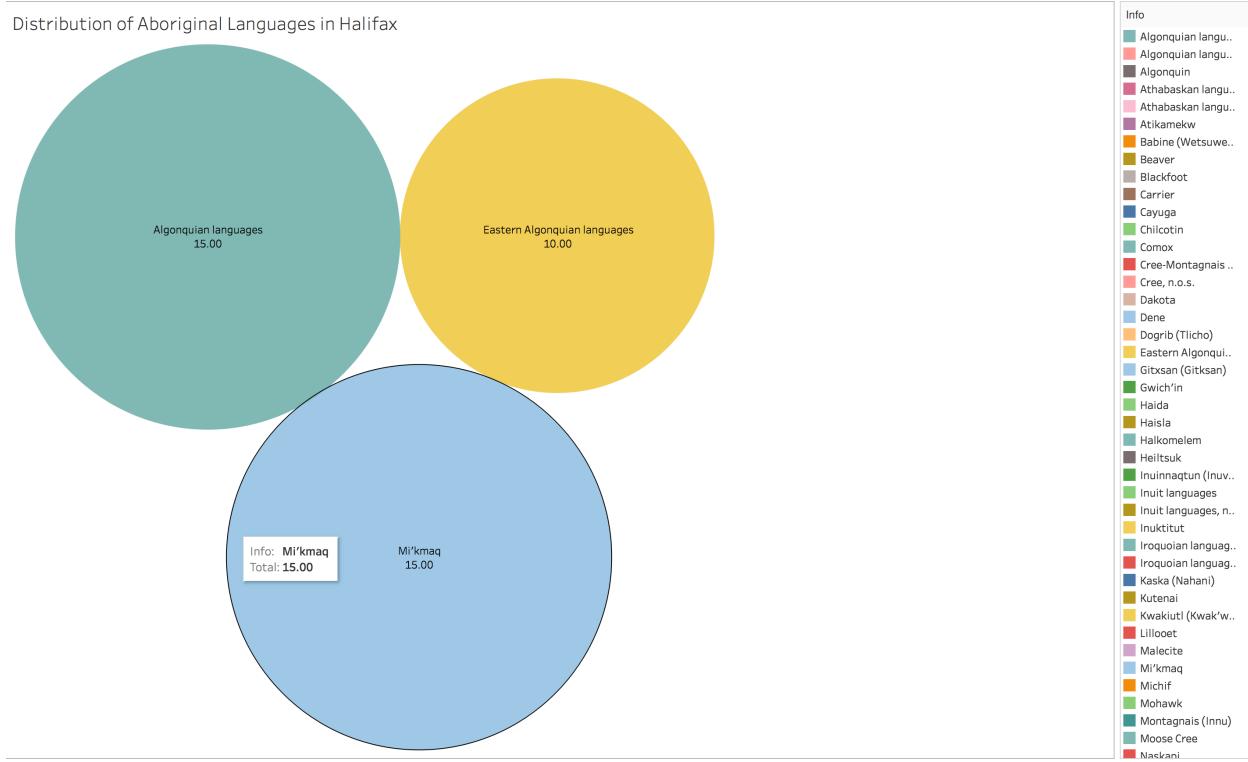
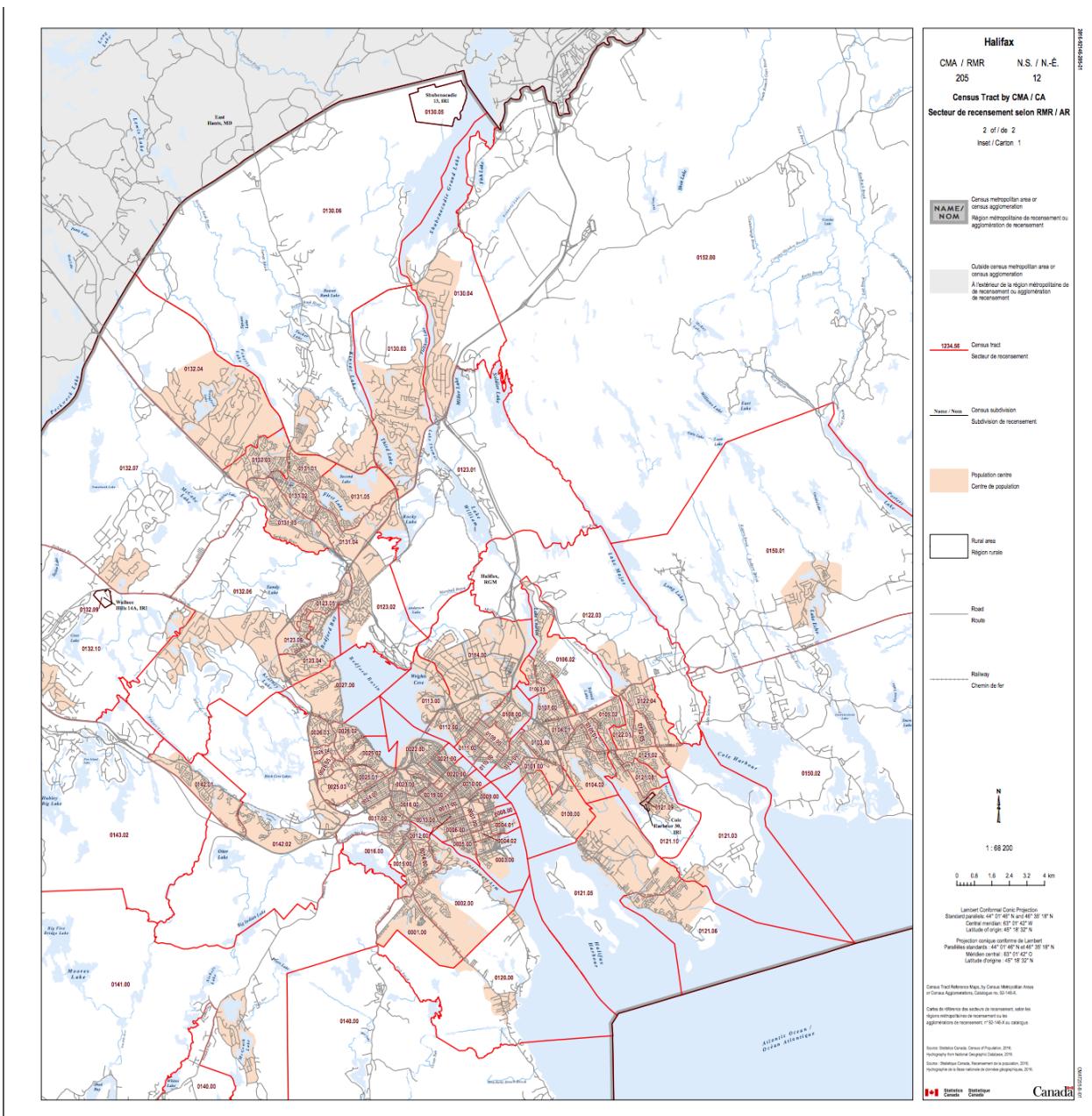


Figure 1: Work Structure Breakdown

As we can see from the above visualizations that Halifax is the most optimal location for a business in the hole of Nova Scotia. We decided to analyze the location within Halifax in more detail. We extracted another dataset specific to regions within Halifax. The below map shows the area division and the associated regional codes below.



The visualization below shows the number of refugees that have move into specific regions of Halifax over the time period between 1980 – 2016.

## Refugee immigration as per census track from 1980 - 2016

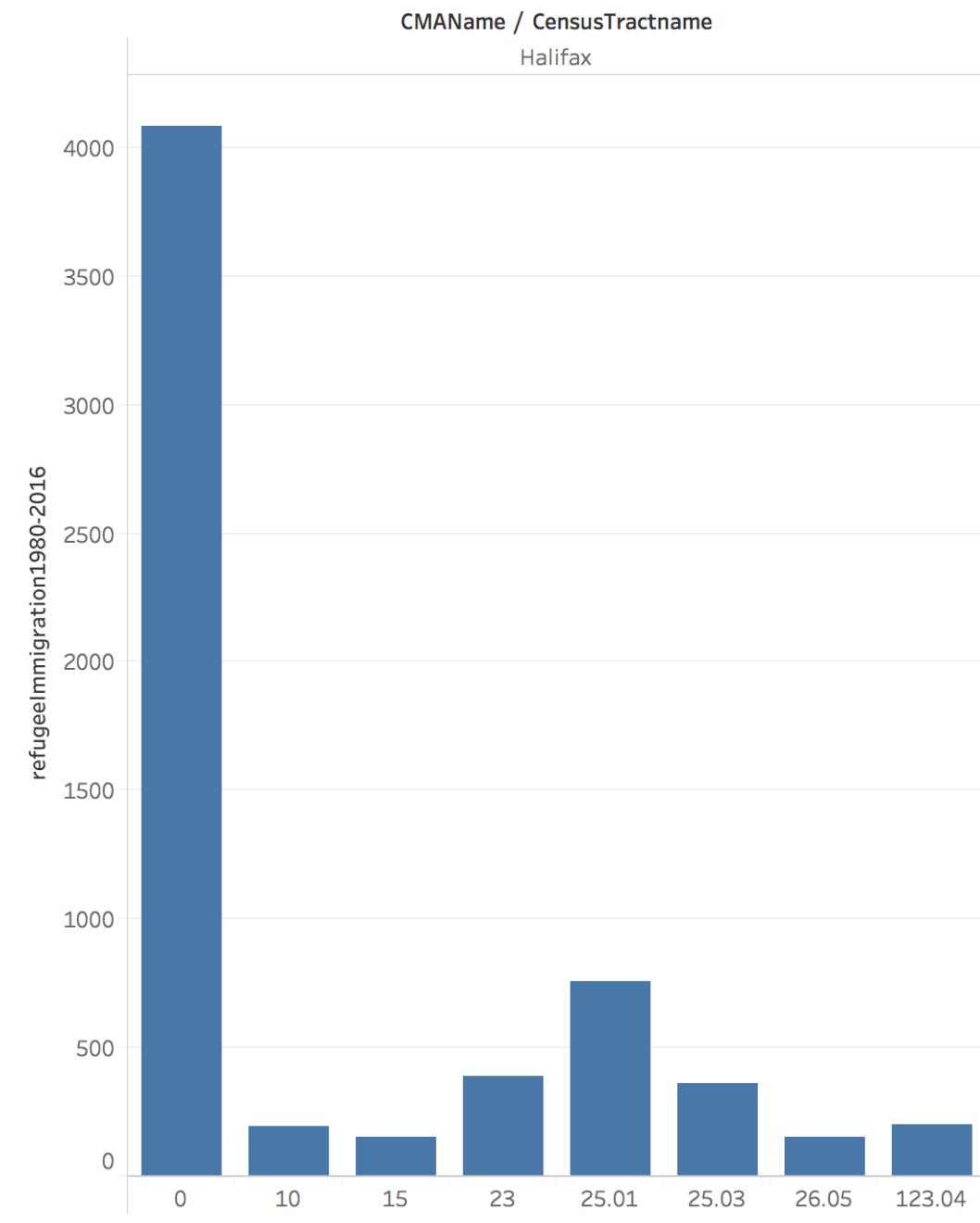


Figure 1: Work Structure Breakdown

These track numbers correspond to regions within Halifax. The zoomed map below shows one of the regions displayed in the bar graph above.

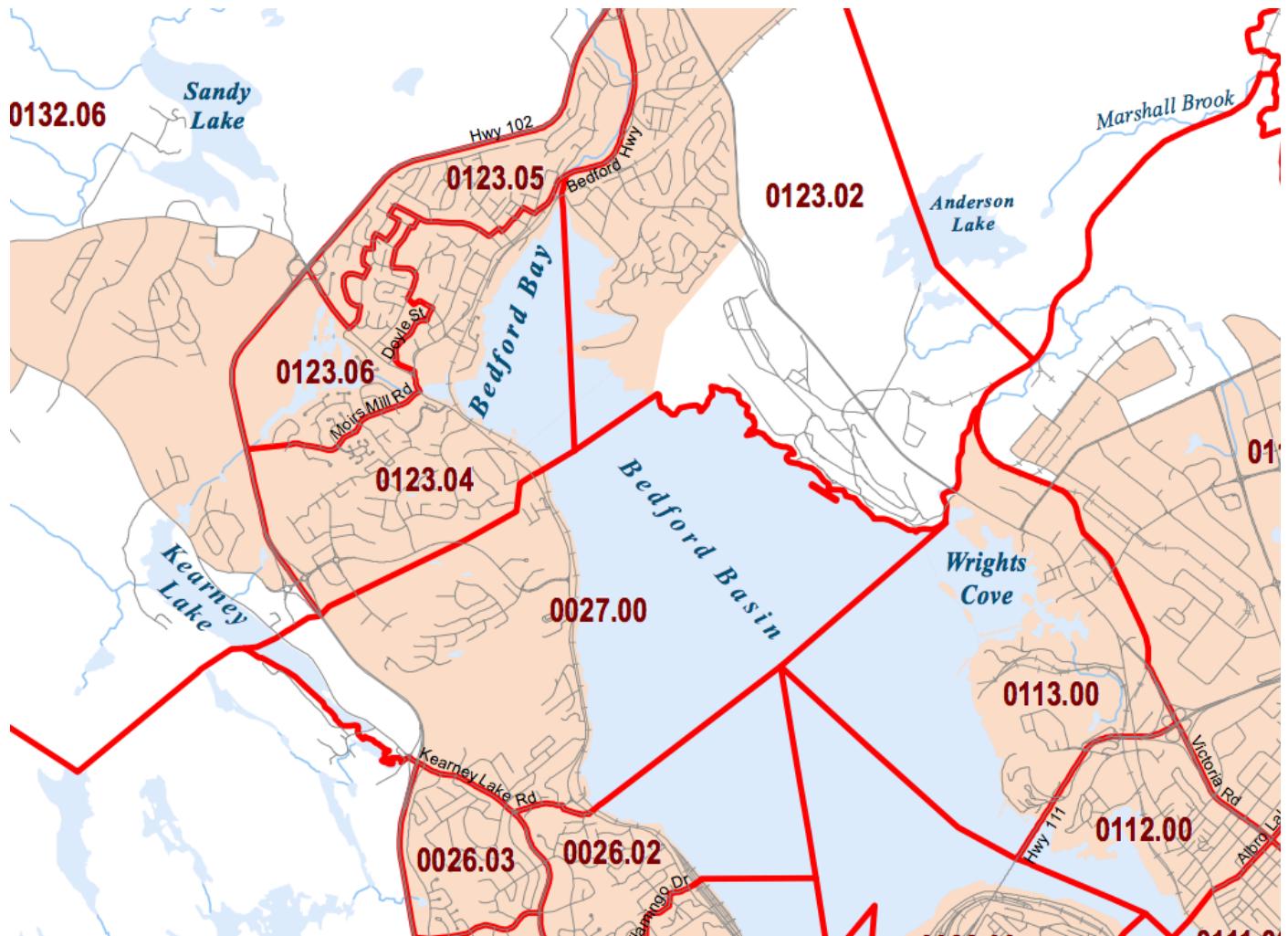


Figure 1: Work Structure Breakdown

Query 3: Find the city with the highest population of senior citizens and a population of people with income higher than \$80,000

```

1 select Geoname, o24_65_years_and_over,aa720_80000_and_over
2 from dwproj.pivot
3 group by o24_65_years_and_over,aa720_80000_and_over DESC;

```

Geoname	o24_65_years_and_over	aa720_80000_and_over
Still Water Lake	90	125
English Corner	145	105
Brookside	180	80
Haves Subdivision	180	35
Port Williams	195	75
Howe's Cove	240	55
Inverness	330	35
Shelburne	335	30
Lake Echo	345	60
Hantsport	370	45
Middleton	400	30
Chester	575	85
Savannah	640	35
Dobie	640	30
Bedford	705	35
Lunenburg	745	70
Pictou	765	70
Bennick	835	30
Enfield - Lantz	880	235
Wolfville	1100	240
Shubenacadie	1165	260
Antigonish	1280	210
New Waterford	1280	165
Yarmouth	1840	155
Amherst	7784	140

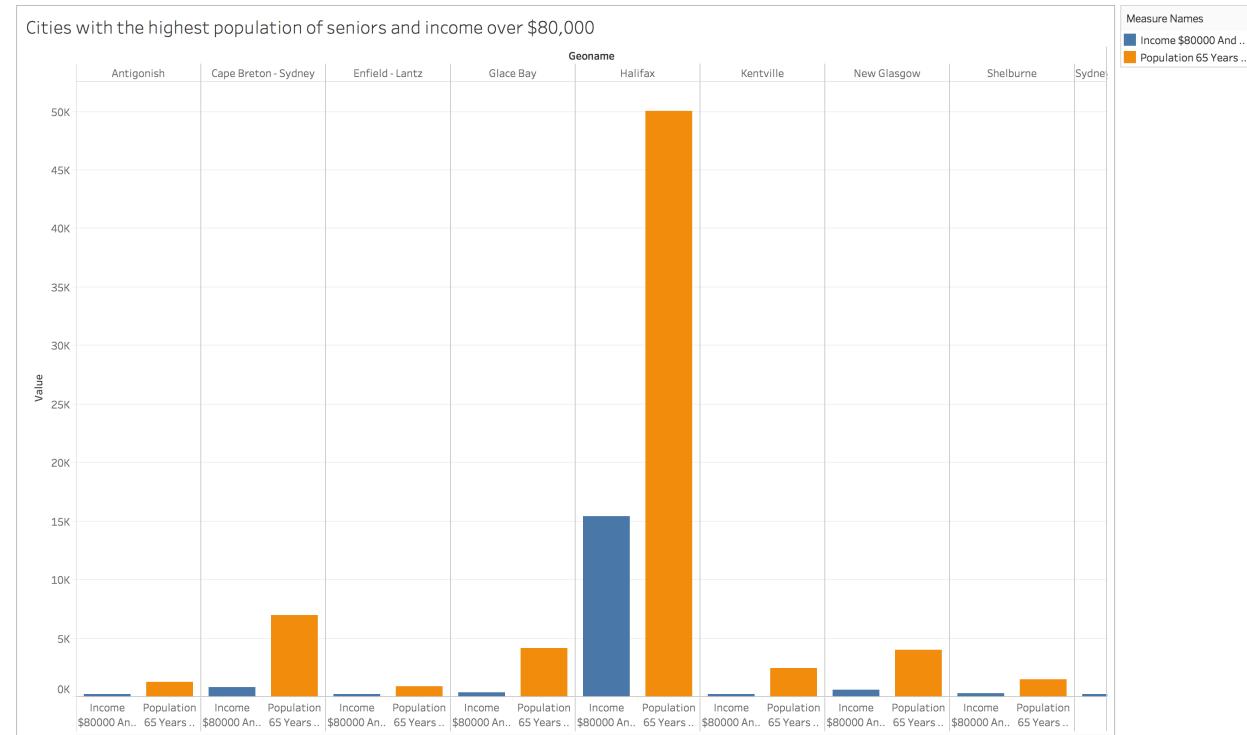


Figure 1: Work Structure Breakdown

This visualization is important to understand which city has the highest growth potential for the senior café. This shows that Halifax has the highest population of both seniors and high income population.

#### 4. Role based work and distribution

Name	Banner ID	Role
Bhavya Chandrappa	B00781097	Data Scientist
Sogra Bilal Memon	B00786252	Data Engineer

Table 1: Team member roles

The roles taken up by Bhavya and Sogra are Data Scientist and Data Engineer respectively. The data was extracted, transformed and loaded to the MySQL Workbench by the Data Scientist. The Data Scientist then queried the results of the resulted dataset. The Data Engineer then took the project forward by analyzing the queries and building graphs to help the customers to look into their searches in a creative view. The Data Engineer also tested the whole process of the project for the user to provide a good user experience.

#### 5. Milestones/Sprints

	Planned	Executed	Tasks
<b>Milestone 1</b>	July 4 - July 7	July 4 July 10	Data gathering and Loading data into EC2 instance
<b>Milestone 2</b>	July 8 - July 14	July 11 – July 20	ETL Process
<b>Milestone 3</b>	July 15 - July 21	July 21 – July 26	Spark SQL query

<b>Milestone 4</b>	July 22 - July 28	July 27 – Aug 1	Tableau Visualization
--------------------	-------------------	-----------------	-----------------------

Table 2: Milestones and sprints

## 6. Limitations of Current work

The limitation of current work includes the fact that the data source is static. Furthermore, the CENSUS data is collected and made available every 5 years this makes it hard to make the dataset dynamic. Another limitation is the memory space. This issue arises because the dataset records have been recorded for the entire Canada. For this pro

## 7. Future Work

For the future work of this project we would like to build a front end where users can select the columns they would like to visualize based on the requirement of their business and the user interface would generate the resulting query and a visualization automatically. In this way a normal user can have access to visualizations that can be shown as part of business pitches and decisions being made without having to know the details about ETL processes or about Data Science.

## 8. Github link

The github link: <https://github.com/SograMemon/DataWarehouse-A2.git>

## 9. Critical Review

The project was designed to help newcomers to find a location to start up their business. The project is quite straightforward as the data is in the same manner. We couldn't experiment with the data. If only there were more required data in the census file, the data generated could have been more creative. We made lots of trials and errors on how to process the data, in the end we were convenient to visualize our results as it made us to improvise for each visualization onwards. We wanted to display the geolocations on the map for the final result as the Geolocations data weren't available in the census file, we stuck to what we have.

## 10. References

- [1] "Tableau Spark SQL Connector Demo - YouTube." [Online]. Available: <https://www.youtube.com/watch?v=OKclf6UdK7c>. [Accessed: 03-Jul-2018].
  - [2] "Spark SQL & DataFrames | Apache Spark." [Online]. Available: <https://spark.apache.org/sql/>. [Accessed: 03-Jul-2018].
  - [3] "90% Of Startups Fail: Here's What You Need To Know About The 10%." [Online]. Available: <https://www.forbes.com/sites/neilpatel/2015/01/16/90-of-startups-will-fail-heres-what-you-need-to-know-about-the-10/#5b657bfa6679>. [Accessed: 03-Jul-2018].
  - [4] "Tableau Help | Tableau Software." [Online]. Available: <https://www.tableau.com/support/help>. [Accessed: 03-Jul-2018].
  - [5] "Answer questions as fast as you can think of them with Tableau | Tableau Software." [Online]. Available: [https://www.tableau.com/trial/tableau-software?utm\\_campaign\\_id=2017049&utm\\_campaign=Prospecting-CORE-ALL-ALL-ALL-ALL&utm\\_medium=Paid+Search&utm\\_source=Google+Search&utm\\_language=EN&utm\\_country=USCA&kw=tableau&adgroup=CTX-Brand-Core-E&adused=221441171506&m](https://www.tableau.com/trial/tableau-software?utm_campaign_id=2017049&utm_campaign=Prospecting-CORE-ALL-ALL-ALL-ALL&utm_medium=Paid+Search&utm_source=Google+Search&utm_language=EN&utm_country=USCA&kw=tableau&adgroup=CTX-Brand-Core-E&adused=221441171506&m). [Accessed: 03-Jul-2018].
  - [6] "Date Functions." [Online]. Available: [https://onlinehelp.tableau.com/current/pro/desktop/en-us/functions\\_functions\\_date.html](https://onlinehelp.tableau.com/current/pro/desktop/en-us/functions_functions_date.html). [Accessed: 03-Jul-2018].
  - [7] "Descriptive Analytics with Tableau - Creating Correlation Matrix and Boxplot Chart - YouTube." [Online]. Available: [https://www.youtube.com/watch?v=eQnH\\_lqUfEA](https://www.youtube.com/watch?v=eQnH_lqUfEA). [Accessed: 03-Jul-2018].
  - [8] Shanelynn, "The Pandas DataFrame – loading, editing, and viewing data in Python," *Shane Lynn*, 18-Dec-2017. [Online]. Available: <https://www.shanelynn.ie/using-pandas-dataframe-creating-editing-viewing-data-in-python/>. [Accessed: 06-Aug-2018].
- "Fixing Column Names in pandas," *The New, Interactive Singles Map.* [Online]. Available: <http://jonathansoma.com/lede/foundations/classes/pandas columns and functions/fixing-column-names-in-pandas/>. [Accessed: 06-Aug-2018].

- “Multiple Criteria Filtering,” *ritchieng.github.io*. [Online]. Available: <https://www.ritchieng.com/pandas-multi-criteria-filtering/>. [Accessed: 06-Aug-2018].
- [9] “Pandas read\_csv: AttributeError: 'NoneType' object has no attribute 'dtype'," *Stack Overflow*. [Online]. Available: <https://stackoverflow.com/questions/30383478/pandas-read-csv-attributeerror-nonetype-object-has-no-attribute-dtype>. [Accessed: 06-Aug-2018].
- [10] “Frequently Asked Questions (FAQ)” *pandas: powerful Python data analysis toolkit - pandas 0.22.0 documentation*. [Online]. Available: <http://pandas.pydata.org/pandas-docs/stable/gotchas.html#support-for-integer-na>. [Accessed: 06-Aug-2018].
- [11] “Convert Pandas column containing NaNs to dtype ‘int’,” *Stack Overflow*. [Online]. Available: <https://stackoverflow.com/questions/21287624/convert-pandas-column-containing-nans-to-dtype-int>. [Accessed: 06-Aug-2018].
- [12] “Select rows from a DataFrame based on values in a column in pandas,” *Stack Overflow*. [Online]. Available: <https://stackoverflow.com/questions/17071871/select-rows-from-a-dataframe-based-on-values-in-a-column-in-pandas>. [Accessed: 06-Aug-2018].
- [13] “pandas dataframe - select by partial string,” *Stack Overflow*. [Online]. Available: <https://stackoverflow.com/questions/11350770/pandas-dataframe-select-by-partial-string>. [Accessed: 06-Aug-2018].
- [14] “Combine two columns of text in dataframe in pandas/python,” *Stack Overflow*. [Online]. Available: <https://stackoverflow.com/questions/19377969/combine-two-columns-of-text-in-dataframe-in-pandas-python>. [Accessed: 06-Aug-2018].
- [15] “pivot\_table No numeric types to aggregate,” *Stack Overflow*. [Online]. Available: <https://stackoverflow.com/questions/39229005/pivot-table-no-numeric-types-to-aggregate>. [Accessed: 06-Aug-2018].