

Introduction: Business Problem

■ The aim of this project is to find a public schools in Chicago, US placed in a suitable location equipped with a proper commercial establishments. In particular this report will be targeted to people moving from other cities/states/countries to Chicago and interested in sending their children to the right school in Chicago.

Data description

- Based on definition of the problem, the following factors that will influence our decision are:
- finding the geographical location of the schools in Chicago,
- ✓ finding the most common venues surrounding a particular school.
- We will be using the geographical coordinates of Chicago and geographical location of the schools to plot school location, and finally cluster our schools 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 City of Chicago containing information on Chicago public schools in 2011-2012 school year, updated in 2018. A dataset consisting of location of the school, its type and other optional parameters describing a school.
- ✓ Part 2: Foursquare API Data
- ✓ Part 3: Creating a new consolidated dataset of the schools, the most common venues and the respective Community Areas along with co-ordinates.: This data will be fetched using Four Square API to explore the venues around schools and to apply machine learning algorithm to cluster the schools and present the findings by plotting it on maps using Folium.

Data description

Using a real-world data set from City of Chicago containing information on Chicago public schools in 2011-2012 school year, updated in 2018.

Chicago Public Schools - Progress Report Cards (2011-2012)

This is a very detailed dataset containing many useful information about each public school in Chicago Some properties of dataset include:

- Name of School
- ▼ Type of School (Elementary, Middle, or High School)
- ✓ Street Address
- ✓ ZIP Code
- ✓ Phone Number
- ✓ Website URL
- ✓ Safety Score
- ✓ Family Involvement Score
- Environment Score
- ✓ Leaders Score
- Teachers Score
- ✓ Latitude
- ✓ Longitude
- ✓ Community Area Name

Data set URL:

https://data.cityofchicago.org/Education/Chicago-Public-Schools-Progress-Report-Cards-2011-/9xs2-f89t

Data description: Foursquare API Data

We will need data about different venues surrounding schools. In order to gain that information, we will use Foursquare locational information. *Foursquare* is a location data provider with information about all manner of venues and events within an area of interest. Such information includes venue names, locations, menus and even photos. As such, the foursquare location platform will be used as the sole data source since all the stated required information can be obtained through the API.

After finding the list of schools, we then connect to the Foursquare API to gather information about venues around every considered school. For each school surrounding, we have chosen the radius to be **100 meters**.

The data retrieved from Foursquare contained information of venues within a specified distance of the longitude and latitude of the school location. The information obtained per venue as follows:

- ✓ Name of School,
- ✓ School Latitude,
- ✓ School Longitude,
- ✓ Name of Venue.
- ✓ Venue Latitude.
- ✓ Venue Longitude,
- ✓ Venue Category.

Based on all the above-described information we have collected a sufficient data to build our model. We cluster the schools together based on similar venue categories. We then present our observations and findings. Using this data, our stakeholders can take the necessary decision.

Methodology

We will be creating our model with the help of Python so we start off by importing all the required packages.

- pandas: To collect and manipulate data in JSON and HTML and then data analysis,
- ✓ Nominatim: To convert an address into latitude and longitude values,
- ✓ folium: Generating maps Chicago,
- ✓ requests: Handle http requests,
- ✓ json_normalize: To normalise data returned by FourSquare API,
- ✓ sklearn: To import Kmeans which is the machine learning model that we are using.

The approach taken here is to explore each surrounding of the considered school, plot the map to show the schools being considered and then build our model by clustering all the similar school surroundings together and finally plot the new map with the clustered schools. We draw insights and then compare and discuss our findings.



Data downloaded from the website of City of Chicago were stored in an csv file:

Reading from dataset:

chicago_schools_df = pd.read_csv('Chicago_Public_Schools_2011-2012_updated.csv', index_col=None, error_bad_lines=False)
chicago_schools_df.head()

	School ID	Name of School	Elementary, Middle, or High School	Street Address	City	State	ZIP Code	Phone Number	Link	Network Manager	 RCDT\$ Code
0	609966	Charles G Hammond Elementary School	ES	2819 W 21st PI	Chicago	IL	60623	(773) 535- 4580	http://schoolreports.cps.edu/SchoolProgressRep	Pilsen-Little Village Elementary Network	 15000000000000000
1	610539	Marvin Camras Elementary School	ES	3000 N Mango Ave	Chicago	IL	60634	(773) 534- 2960	http://schoolreports.cps.edu/SchoolProgressRep	Fullerton Elementary Network	 15000000000000000
2	609852	Eliza Chappell Elementary School	ES	2135 W Foster Ave	Chicago	IL	60625	(773) 534- 2390	http://schoolreports.cps.edu/SchoolProgressRep	Ravenswood- Ridge Elementary Network	 15000000000000000
3	609835	Daniel R Cameron Elementary School	ES	1234 N Monticello Ave	Chicago	IL	60651	(773) 534- 4290	http://schoolreports.cps.edu/SchoolProgressRep	Garfield- Humboldt Elementary Network	 1500000000000000
4	610521	Sir Miles Davis Magnet Elementary Academy	ES	6730 S Paulina St	Chicago	IL	60636	(773) 535- 9120	http://schoolreports.cps.edu/SchoolProgressRep	Englewood- Gresham Elementary Network	 1500000000000000

5 rows x 79 columns

.



This is a very detailed dataset, and we extracted the information which is the most useful to us:

	Name of School	Street Address	City	ZIP Code	Elementary, Middle, or High School	Safety Score	Latitude	Longitude	Community Area Name
0	Charles G Hammond Elementary School	2819 W 21st PI	Chicago	60623	ES	40.0	41.852691	-87.696278	SOUTH LAWNDALE
1	Marvin Camras Elementary School	3000 N Mango Ave	Chicago	60634	ES	54.0	41.934966	-87.770165	BELMONT CRAGIN
2	Eliza Chappell Elementary School	2135 W Foster Ave	Chicago	60625	ES	70.0	41.975867	-87.683254	LINCOLN SQUARE
3	Daniel R Cameron Elementary School	1234 N Monticello Ave	Chicago	60651	ES	42.0	41.903785	-87.717963	HUMBOLDT PARK
4	Sir Miles Davis Magnet Elementary Academy	6730 S Paulina St	Chicago	60636	ES	35.0	41.771222	-87.666567	WEST ENGLEWOOD

We start data exploration by finding how many schools exists in

How many schools in each Chicago Community Area:

```
chicago_schools_sel['Community Area Name'].value_counts()

AUSTIN 23
SOUTH LAWNDALE 22
WEST TOWN 20
ENGLEWOOD 17
NEAR WEST SIDE 16
...
BURNSIDE 1
MONTCLARE 1
LOOP 1
OAKLAND 1
OHARE 1
Name: Community Area Name, Length: 77, dtype: int64
```

There are three types of public schools in Chicago: Elementary (ES), Middle (MS) and High (HS) schools:

How many public schools of a particular type are in Chicago:

```
seriesObjE = chicago_schools_sel.apply(lambda x: True if x['Elementary, Middle, or High School'] == 'ES' else False , axis=1)
# Count number of True in series
numOfRowsE = len(seriesObjE[seriesObjE == True].index)
print('Number of elementary schools : ', numOfRowsE)

seriesObjM = chicago_schools_sel.apply(lambda x: True if x['Elementary, Middle, or High School'] == 'MS' else False , axis=1)
numOfRowsM = len(seriesObjM[seriesObjM == True].index)
print('Number of middle schools : ', numOfRowsM)

seriesObjH = chicago_schools_sel.apply(lambda x: True if x['Elementary, Middle, or High School'] == 'HS' else False , axis=1)
numOfRowsH = len(seriesObjH[seriesObjH == True].index)
print('Number of high schools : ', numOfRowsH)

Number of elementary schools : 462
Number of middle schools : 11
Number of high schools : 93
```

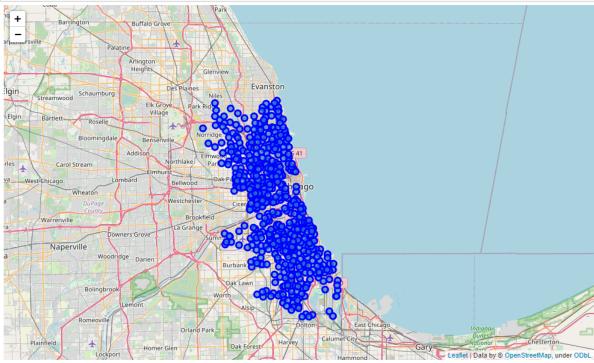
Next, we use *Nomatim geolocator* to find geographical coordinates of Chicago, what will be needed to plot map of Chicago together with school's locations overlaid on it.

Using geopy Nominatim geolocator to fing geographical coordinates of Chicago:

```
address = "Chicago, IL"

geolocator = Nominatim(user_agent="chicago_explorer")
location = geolocator.geocode(address)
latitude = location.latitude
longitude = location.longitude
print('The geograpical coordinates of Chicago are {}, {}.'.format(latitude, longitude))
```

The geograpical coordinates of Chicago are 41.8755616, -87.6244212.



Due to substantial number of public schools in Chicago, we focus on *high schools* only.

```
: for lat, lng, ncomm, nschool in zip(
           chicago_high_schools['Latitude'],
chicago_high_schools['Longitude'],
chicago_high_schools['Community Area Name'],
           chicago_high_schools['Name of School']):
       label = '{}, {}'.format(ncomm, nschool)
       label = folium.Popup(label, parse_html=True)
       folium.CircleMarker(
           [lat, lng],
           radius=5,
           popup=label,
           color='blue',
           fill=True,
           fill_color='#3186cc',
           fill_opacity=0.7,
           parse_html=False).add_to(map_chicago_high)
  map_chicago_high
               Naperville
```

Next, we prepare a function which uses *Foursquare API* to fetch venues around a given location and use it.

Fetching venues around public high schools:

```
chicago school venues = getNearbyVenues(names=chicago high schools['Name of School'],
                                   latitudes=chicago_high_schools['Latitude'],
                                   longitudes=chicago high schools['Longitude']
Walter Payton College Preparatory High School
Manley Career Academy High School
Northside College Preparatory High School
Michele Clark Academic Prep Magnet High School
Uplift Community High School
Morgan Park High School
Bronzeville Scholastic Academy High School
William J Bogan High School
Emil G Hirsch Metropolitan High School
Austin Polytechnical Academy High School
World Language Academy High School
Multicultural Academy of Scholarship
Mason High School
Marie Sklodowska Curie Metropolitan High School
George Washington High School
Robert Lindblom Math & Science Academy High School
Benito Juarez Community Academy High School
Hyde Park Academy High School
John Marshall Metropolitan High School
Friedrich W von Steuben Metropolitan Science High School
Southside Occupational Academy High School
Chicago Military Academy High School
Eric Solorio Academy High School
Neal F Simeon Career Academy High School
John Hancock College Preparatory High School
Roald Amundsen High School
Edwin G Foreman High School
Paul Laurence Dunbar Career Academy High School
Charles P Steinmetz Academic Centre High School
Gurdon S Hubbard High School
Albert G Lane Technical High School
Carl Schurz High School
Dyett High School
Phoenix Military Academy High School
Chicago Vocational Career Academy High School
```

The collected data we put into a new data frame and then we group them with respect to school location.

Groupping of venues with respect to school location: chicago_school_venues.groupby('Name of School').count().drop(['School Latitude', 'School Longitude', 'Venue Category'], axis = 1) Alcott High School for the Humanities Benito Juarez Community Academy High School Carl Schurz High School Chicago High School for Agricultural Sciences Chicago Military Academy High School DeVry University Advantage Academy High School Friedrich W von Steuben Metropolitan Science High School Gage Park High School George H Corliss High School Gwendolyn Brooks College Preparatory Academy High School Hyman G Rickover Naval Academy High School Lincoln Park High School Marie Sklodowska Curie Metropolitan High School Michele Clark Academic Prep Magnet High School New Millennium High School of Health at Bowen Nicholas Senn High School Northside College Preparatory High School Orr Academy High School Roald Amundsen High School Roberto Clemente Community Academy High School Spry Community Links High School Stephen T Mather High School Thomas Kelly High School Uplift Community High School Wells Community Academy High School William J Bogan High School William Jones College Preparatory High School print('There are {} uniques categories.'.format(len(chicago_school_venues['Venue Category'].unique())))

There are 47 uniques categories.

Since we are trying to find out what are the different kinds of venue categories present in each high school surrounding and then calculate the top 5 common venues to base our similarity on, we use the One Hot Encoding to work with our categorical datatype of the venue categories. This helps to convert the categorical data into numeric data.

Top 5 most common venues around high school:

```
num_top_venues = 5
for school in chicago_school_grouped['Name of School']:
   print("----"+school+"----")
   temp = chicago school grouped[chicago school grouped['Name of School'] == school].T.reset index()
   temp.columns = ['venue', 'freq']
   temp = temp.iloc[1:]
   temp['freq'] = temp['freq'].astype(float)
   temp = temp.round({'freq': 2})
   print(temp.sort values('freq', ascending=False).reset index(drop=True).head(num top venues))
   print('\n')
----Alcott High School for the Humanities----
              Dog Run 0.5
     Korean Restaurant 0.0
----Benito Juarez Community Academy High School----
   Mexican Restaurant
   American Restaurant
           Pizza Place
     Korean Restaurant
               Market 0.0
----Carl Schurz High School----
```

We will be using KMeans Clustering Machine learning algorithm to cluster similar school's surroundings together. We will be going with the number of clusters as 5.

Clustering Chicago high schools:

We will be using KMeans Clustering Machine learning algorithm to cluster similar school's surroundings together. We will be going with the number of clusters as 5.

```
# create map
map clusters = folium.Map(location=[latitude, longitude], zoom start=12)
# set color scheme for the clusters
x = np.arange(kclusters)
ys = [i + x + (i*x)**2 \text{ for } i \text{ in } range(kclusters)]
colors_array = cm.rainbow(np.linspace(0, 1, len(ys)))
rainbow = [colors.rgb2hex(i) for i in colors array]
# add markers to the map
markers_colors = []
for lat, lon, poi, cluster in zip(chicago_merged['Latitude'], chicago_merged['Longitude'], chicago_merged['Name of School'], chicago_merged['Name of Name of School'], chicago_merged['Name of Name 
             label = folium.Popup(str(poi) + ' Cluster ' + str(cluster), parse html=True)
              folium.CircleMarker(
                             [lat, lon],
                             radius=5,
                             popup=label,
                             color=rainbow[cluster-1],
                             fill_color=rainbow[cluster-1],
                            fill_opacity=0.7).add_to(map_clusters)
map_clusters
                                                                      West Garfield
```

Cluster analysis:

Cluster 1:

chicago_merged.loc[chicago_merged['Cluster Labels'] == 0, chicago_merged.columns[[1] + list(range(5, chicago_merged.shape[1]))]]

	Street Address	Safety Score	Latitude	Longitude	Community Area Name	Cluster Labels	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue	4th Most Common Venue	5th Most Common Venue
17	5501 N Kedzie Ave	99.0	41.981352	-87.708672	NORTH PARK	0	Music Venue	Women's Store	Historic Site	Furniture / Home Store	Fried Chicken Joint
28	5101 W Harrison St	NaN	41.872857	-87.753355	AUSTIN	0	Gym / Fitness Center	Coffee Shop	Furniture / Home Store	Fried Chicken Joint	Football Stadium
30	900 W Wilson Ave	50.0	41.965574	-87.652522	UPTOWN	0	Pharmacy	Café	Coffee Shop	Furniture / Home Store	Fried Chicken Joint
38	3939 W 79th St	20.0	41.749348	-87.721097	ASHBURN	0	Furniture / Home Store	Fast Food Restaurant	Discount Store	Women's Store	Coffee Shop
77	2150 S Laflin St	46.0	41.852673	-87.663769	LOWER WEST SIDE	0	Mexican Restaurant	Women's Store	Historic Site	Furniture / Home Store	Fried Chicken Joint
95	5039 N Kimball Ave	70.0	41.973193	-87.713350	NORTH PARK	0	River	Bus Station	Women's Store	Coffee Shop	Furniture / Home Store
107	3519 S Giles Ave	32.0	41.830538	-87.619178	DOUGLAS	0	History Museum	Pizza Place	Historic Site	Wings Joint	Cosmetics Shop
116	5110 N Damen Ave	51.0	41.975079	-87.679521	LINCOLN SQUARE	0	Basketball Court	Pool	Women's Store	Coffee Shop	Furniture / Home Store
139	3601 N Milwaukee Ave	48.0	41.946408	-87.735625	IRVING PARK	0	Thai Restaurant	Asian Restaurant	Martial Arts School	Convenience Store	Women's Store
190	4015 N Ashland Ave	64.0	41.954784	-87.668916	LAKE VIEW	0	Chinese Restaurant	Coffee Shop	Thai Restaurant	Fried Chicken Joint	Breakfast Spot
238	3300 N Campbell	NaN	41.941426	-87.690799	NORTH CENTER	0	Furniture / Home Store	Salon / Barbershop	Women's Store	Coffee Shop	Fried Chicken Joint
264	5630 S Rockwell St	14.0	41.791014	-87.688991	GAGE PARK	0	Clothing Store	Coffee Shop	Furniture / Home Store	Fried Chicken Joint	Football Stadium
298	2001 N Orchard St	65.0	41.918304	-87.645974	LINCOLN PARK	0	Women's Store	Art Gallery	BBQ Joint	Burger Joint	Mediterranean Restaurant
311	730 N Pulaski Rd	NaN	41.894448	-87.726203	HUMBOLDT PARK	0	Fast Food Restaurant	Women's Store	Coffee Shop	Furniture / Home Store	Fried Chicken Joint

Cluster analysis:

Cluster 2:

chicago_merged.loc[chicago_merged['Cluster Labels'] == 1, chicago_merged.columns[[1] + list(range(5, chicago_merged.shape[1]))]]

	Street Address	Safety Score	Latitude	Longitude	Community Area Name	Cluster Labels	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue	4th Most Common Venue	5th Most Common Venue
244	4136 S California Ave	36.0	41.818711	-87.694675	BRIGHTON PARK	1	Park	Women's Store	Historic Site	Furniture / Home Store	Fried Chicken Joint
463	2957 N Hoyne Ave	70.0	41.935761	-87.680524	NORTH CENTER	1	Dog Run	Park	Women's Store	Historic Site	Furniture / Home Store
473	250 E 111th St	64.0	41.692790	-87.616381	ROSELAND	1	Park	Women's Store	Historic Site	Furniture / Home Store	Fried Chicken Joint

Cluster 3:

chicago_merged.loc[chicago_merged['Cluster Labels'] == 2, chicago_merged.columns[[1] + list(range(5, chicago_merged.shape[1]))]]

	Street Address	Safety Score	Latitude	Longitude	Community Area Name	Cluster Labels	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue	4th Most Common Venue	5th Most Common Venue
198	821 E 103rd St	33.0	41.707391	-87.603078	PULLMAN	2	Football Stadium	Women's Store	Coffee Shop	Furniture / Home Store	Fried Chicken Joint

Cluster analysis:

Cluster 4:

chicago_merged.loc[chicago_merged['Cluster Labels'] == 3, chicago_merged.columns[[1] + list(range(5, chicago_merged.shape[1]))]]

	Street Address	Safety Score	Latitude	Longitude	Community Area Name	Cluster Labels	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue	4th Most Common Venue	5th Most Common Venue
68	4959 S Archer Ave	43.0	41.803046	-87.722007	ARCHER HEIGHTS	3	Hotel	Women's Store	Historic Site	Furniture / Home Store	Fried Chicken Joint
176	5900 N Glenwood Ave	64.0	41.989051	-87.665262	EDGEWATER	3	Hotel	Women's Store	Historic Site	Furniture / Home Store	Fried Chicken Joint
274	5900 N Glenwood Ave	48.0	41.989051	-87.665262	EDGEWATER	3	Hotel	Women's Store	Historic Site	Furniture / Home Store	Fried Chicken Joint

Cluster 5:

chicago_merged.loc[chicago_merged['Cluster Labels'] == 4, chicago_merged.columns[[1] + list(range(5, chicago_merged.shape[1]))]]

	Stre Addres		Latitude	Longitude	Community Area Name	Cluster Labels	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue	4th Most Common Venue	5th Most Common Venue
5	2710 89th	E 17.0	41.733761	-87.557753	SOUTH CHICAGO	4	American Restaurant	Coffee Shop	Furniture / Home Store	Fried Chicken Joint	Football Stadium

Results and Discussion

The object of the business problem was to help Chicago migrants to identify suitable public school to their children, located in area surrounded with the appropriate venues. This has been achieved by first making use of Chicago Public Schools data to identify a proper place with considerable number of venues. Due to substantial number of public schools in Chicago focus was made on the public high schools only. Next, grouping of the high schools into clusters was done to assist the migrants by providing them with relevant data about venues and safety of a given school surrounding.

Conclusion

We have explored the Chicago Public Schools data to understand different types of public schools in all Community Areas of Chicago and later categorized them into different types. This helped us group the schools. We further shortlist the high schools based on the common venues, to choose clusters of schools which best suits the business problem.