

Helsingborg City Data

Data published by the city of Helsingborg is used for this project. It is hard to get a sense of raw data. The data consist of geospatial locations of different business point in different categories like restaurants, cafe, etc from Helsingborg. Preprocessing of data is done using Python libraries. K-Means clustering is performed on this data to see the patterns and the results are visualized using Folium map interface. optimal k value for k-means is chosen using the Elbow method. The data is so complex that it is really hard to see the optimal value from the Elbow method. Later in this project foursquare API is used to fetch data for the same postal codes and clustered data using the above methodology and the results are shown below.

```
import numpy as np
import pandas as pd
import json
import matplotlib as mlt
import matplotlib.pyplot as plt
import matplotlib.cm as cm
import matplotlib.colors as colors
from sklearn.cluster import KMeans
```

Python libraries used in this project

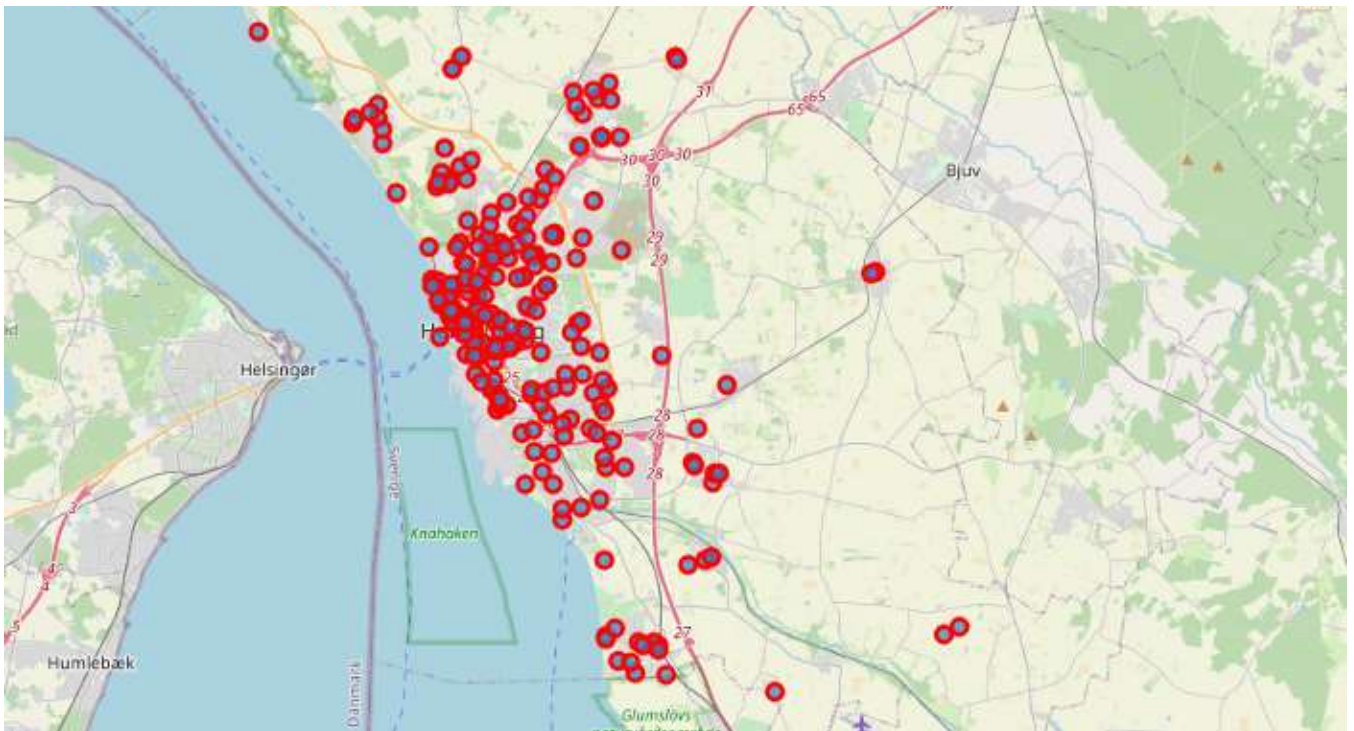
Sample data

	CITY	POSTCODE	geo_shape		geo_point_2d		area
0	Helsingborg	25482	{"type": "Polygon", "coordinates": [[[12.65947...		56.0894565738,12.6568019271		Hittarp - Laröd
1	Helsingborg	25361	{"type": "Polygon", "coordinates": [[[12.74727...		56.0132736825,12.7530666401		Ättekulla

KATEGORI	NAMN	ADRESS	HEMSIDA	Kategori webb	Exportdatum	geo_shape	geo_point_2d	area	postalcode
Buliker	Borgmans Fisk Etr	Kopparmöllegatan 17	http://www.borgmansfisk.net/	Buliker	2020-07-16	{"type": "Point", "coordinates": [12.697555275...	56.0534618304,12.6975552755	Slottshöjden	25435
Buliker	Polshop	Södergatan 16	http://www.polshop.se/	Buliker	2020-07-16	{"type": "Point", "coordinates": [12.701065141...	56.0421044608,12.7010651416	Söder	25225

Data that we are interested in is Geospatial data, Category of each entry, name.





Data before clustering

One hot encoding is performed on categorical variables to feed it to the clustering algorithm

Butiker	Caféer	Personalrestauranger	Restauranger	Skolrestauranger	Snabbmat	Vårdverksamheter
1	0	0	0	0	0	0
1	0	0	0	0	0	0
1	0	0	0	0	0	0
1	0	0	0	0	0	0
1	0	0	0	0	0	0

Data is grouped over postal codes and mean is performed at the grouped level

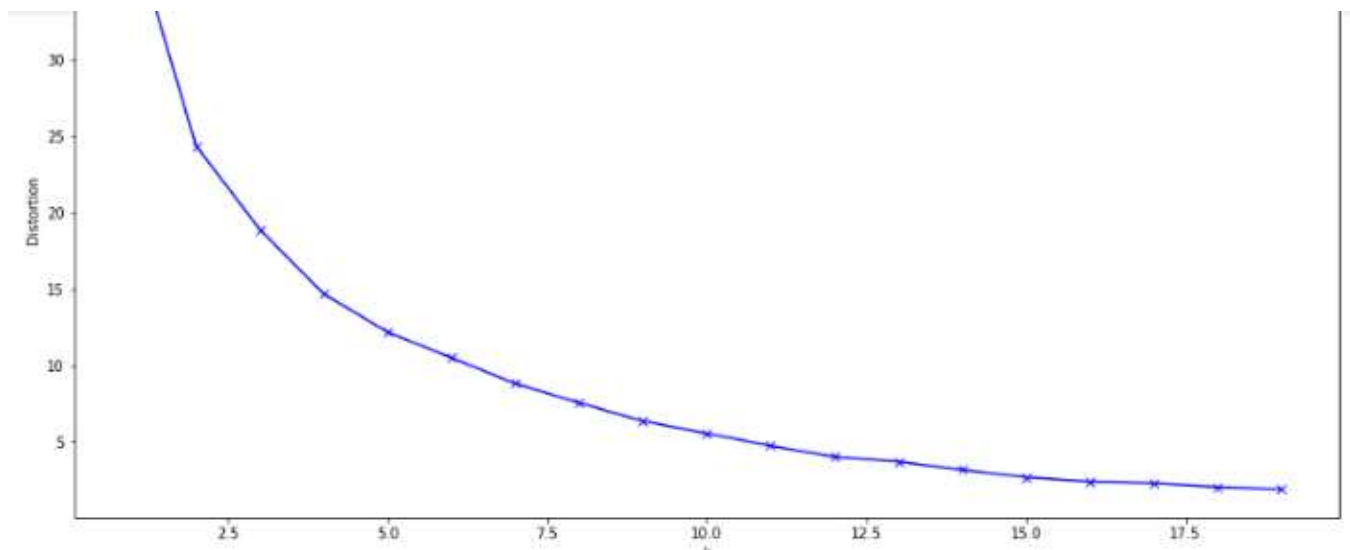
postalcode	Butiker	Caféer	Personalrestauranger	Restauranger	Skolrestauranger	Snabbmat	Vårdverksamheter
25220	0.000000	0.000000	0.000000	1.000000	0.000000	0.000000	0.000000
25221	0.000000	0.250000	0.000000	0.583333	0.083333	0.083333	0.000000
25222	0.000000	1.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25223	0.090909	0.090909	0.090909	0.272727	0.272727	0.181818	0.000000
25224	0.000000	0.000000	0.000000	0.777778	0.000000	0.111111	0.111111

the top places for each postal code are performed

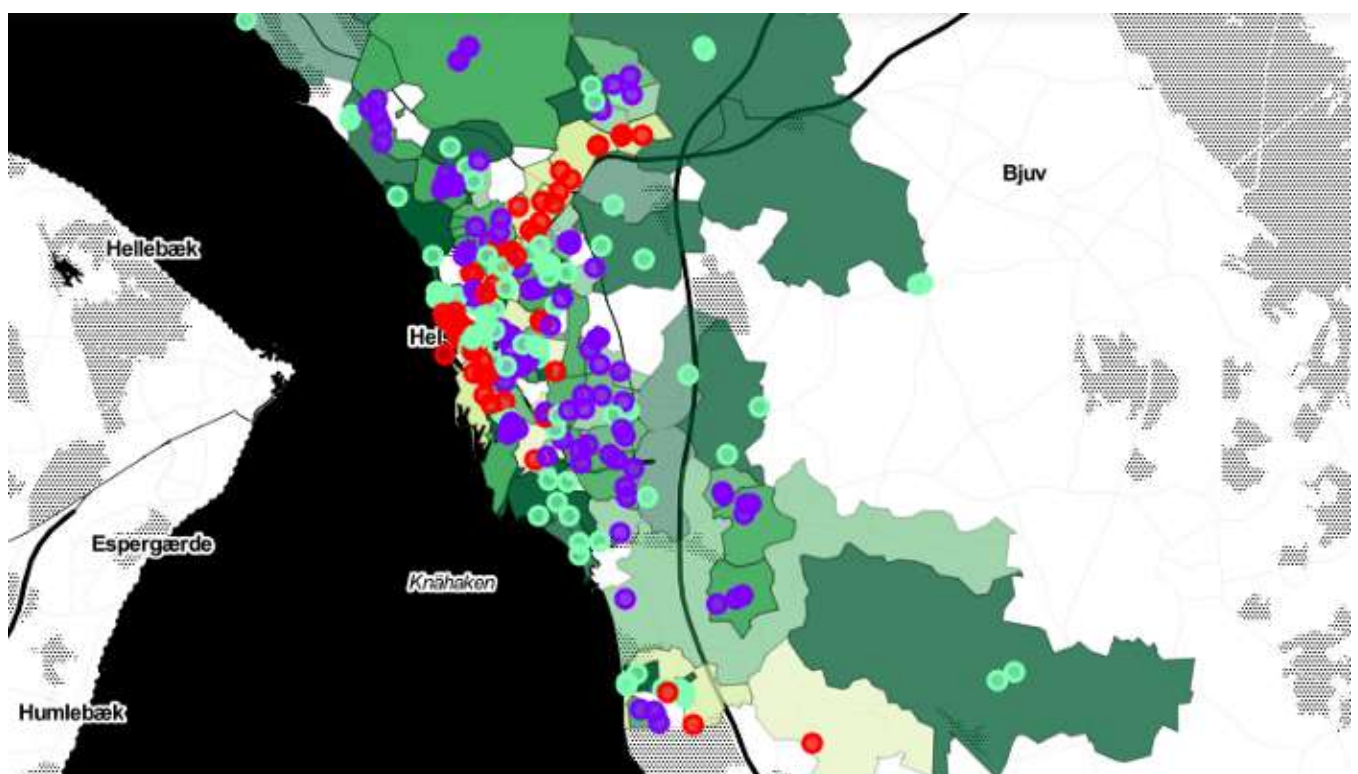
postalcode	1st	2nd	3rd	4th	5th
25220	Restauranger	Vårdverksamheter	Snabbmat	Skolrestauranger	Personalrestauranger
25221	Restauranger	Caféer	Snabbmat	Skolrestauranger	Vårdverksamheter
25222	Caféer	Vårdverksamheter	Snabbmat	Skolrestauranger	Restauranger
25223	Skolrestauranger	Restauranger	Snabbmat	Personalrestauranger	Caféer
25224	Restauranger	Vårdverksamheter	Snabbmat	Skolrestauranger	Personalrestauranger

Find optimal k value by The Elbow Method

It was difficult to choose an optimal value for this dataset



3 Clusters data is visualized by Folium maps



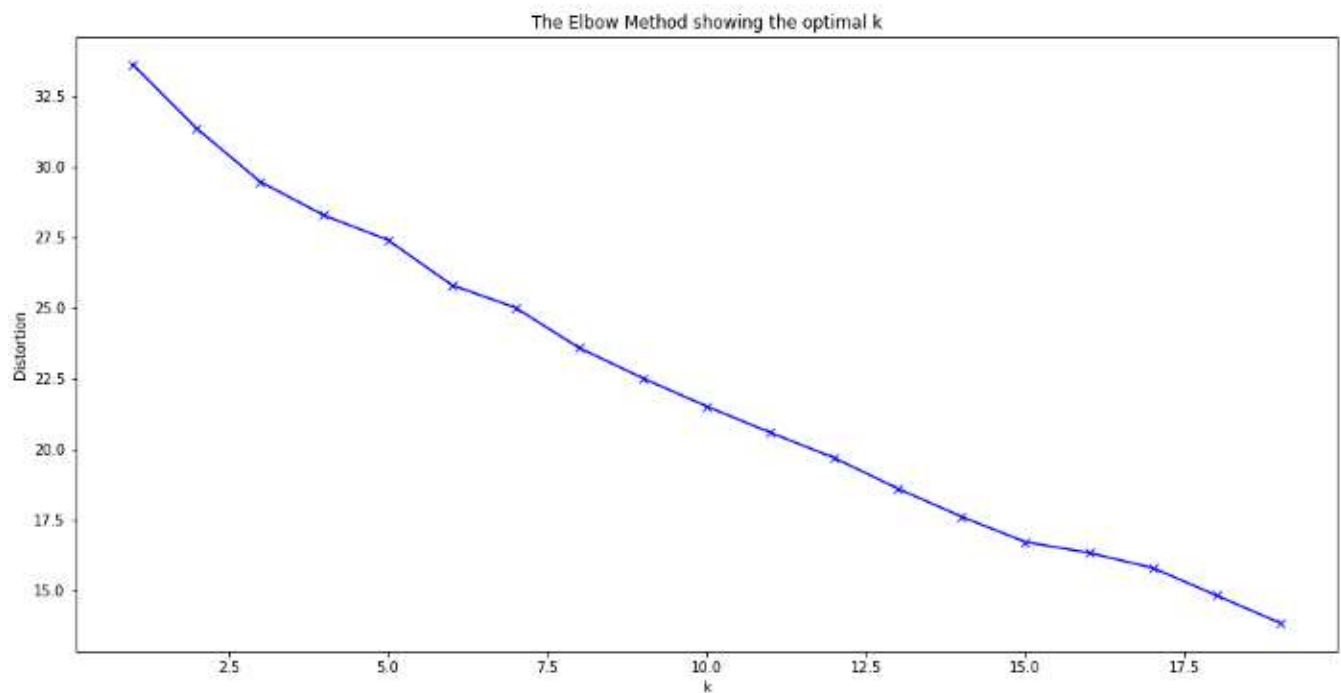
Data grouped on postal code and mean is performed on group-level data.

POSTCODE	American Restaurant	Art Museum	Asian Restaurant	Athletics & Sports	Auto Dealership	Auto Garage	Bagel Shop	Bakery	Bar	Beach	Beach Bar	Beer Bar	Bistro	Bookstore
25220	0.0	0.020833	0.0	0.0	0.0	0.0	0.0	0.020833	0.041667	0.0	0.0	0.020833	0.020833	0.020833
25221	0.0	0.025641	0.0	0.0	0.0	0.0	0.0	0.000000	0.051282	0.0	0.0	0.025641	0.025641	0.025641

Most common business places based on the data from Foursquare API for each postal code of Helsingborg

POSTCODE	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue	4th Most Common Venue	5th Most Common Venue	6th Most Common Venue	7th Most Common Venue	8th Most Common Venue	9th Most Common Venue	10th Most Common Venue
25220	Hotel	Scandinavian Restaurant	Restaurant	Café	Coffee Shop	Burger Joint	Ice Cream Shop	Salad Place	Bar	Italian Restaurant
25221	Scandinavian Restaurant	Restaurant	Pizza Place	Bar	Burger Joint	Hotel	Ice Cream Shop	Salad Place	Nightclub	Café
25222	Park	Scandinavian Restaurant	Italian Restaurant	Grocery Store	Playground	Gym	Deli / Bodega	Pizza Place	Plaza	Coffee Shop
25223	Coffee Shop	Hotel	Café	Restaurant	Italian Restaurant	Park	Asian Restaurant	Pub	Scandinavian Restaurant	Thai Restaurant
25224	Hotel	Restaurant	Coffee Shop	Bar	Café	Thai Restaurant	Scandinavian Restaurant	Asian Restaurant	Pub	Italian Restaurant

picking an optimal k-value was not easy for this dataset

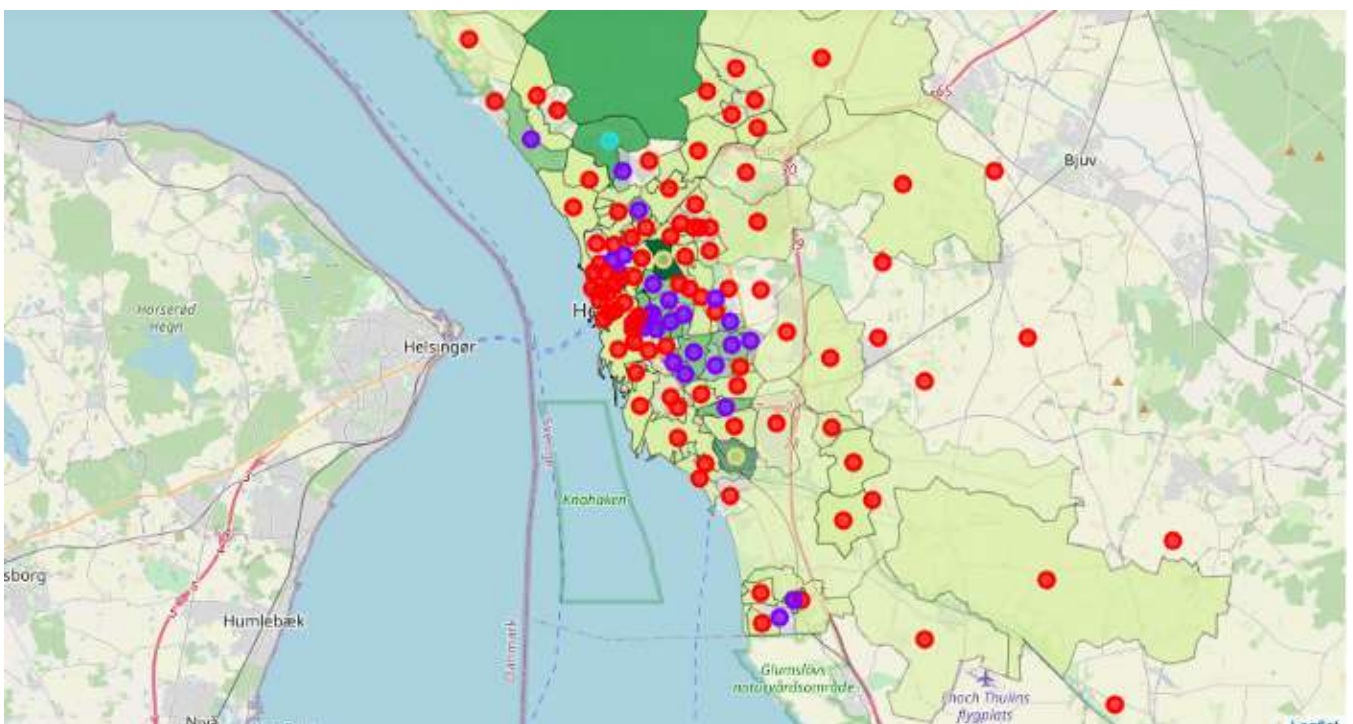


Finally, data with clusters are merged

POSTCODE	geo_shape	geo_point_2d	area	type	coordinates	polygon_coordinates	Cluster Labels	1st Most Common Venue	2nd Most Common Venue
----------	-----------	--------------	------	------	-------------	---------------------	----------------	-----------------------	-----------------------

25482	<pre>{ "type": "Polygon", "coordinates": [[[12.65947...</pre>	56.0894565738,12.6568019271	Hiltarp - Laröd	Polygon	<pre>[[[12.65947745564023, 56.08065244077402], [12....</pre>	<pre>[[12.65947745564023, 56.08065244077402], [12.6...</pre>	1	Playground	Pizza Place
25361	<pre>{ "type": "Polygon", "coordinates": [[[12.74727...</pre>	56.0132736825,12.7530666401	Ättekulla	Polygon	<pre>[[[12.747272091715816, 56.01611001498235], [12....</pre>	<pre>[[12.747272091715816, 56.01611001498235], [12....</pre>	0	Brewery	Supermarket
25465	<pre>{ "type": "Polygon", "coordinates": [[[12.72806...</pre>	56.0721772539,12.7345182361	Dalhem	Polygon	<pre>[[[12.728065859762449, 56.07223268630626], [12....</pre>	<pre>[[12.728065859762449, 56.07223268630626], [12....</pre>	0	Playground	Shopping Mall
25665	<pre>{ "type": "Polygon", "coordinates": [[[12.74776...</pre>	56.0499416462,12.7499930315	Adolfsberg	Polygon	<pre>[[[12.747764472346708, 56.04706467772815], [12....</pre>	<pre>[[12.747764472346708, 56.04706467772815], [12....</pre>	0	Business Service	Athletics & Sports
25244	<pre>{ "type": "MultiPolygon", "coordinates": [[[12....</pre>	56.0413191382,12.7085138805	Eneborg	MultiPolygon	<pre>[[[12.709012002277971, 56.04312887626643], [1....</pre>	<pre>[[12.709012002277971, 56.04312887626643], [12....</pre>	0	Bus Stop	Pizza Place

4 clusters with the Foursquare data



Conclusions:

Data preprocessing is done in Python and prepared data for k-means clustering to discover patterns from complex unlabeled data and the results are visualized using the folium maps. If we see the clustered data from the first experiment, red circle clusters are the areas where the main city centre is located and most of the popular restaurants, supermarkets are located and also most populated areas. Foursquare API data clearly lacks distinguishable clusters, require more data and good features to perform better in discovering patterns.