

Introduction/Business Problem

For my project I decided to explore whether there is a correlation between neighborhood types derived using ML algorithms and different crime types also broken down by some clusters using K-means. To me this looks to be a very interesting thing to explore: maybe there are some hidden dependencies that we can discover. For example: liquor shops may draw more DUI crimes nearby. In ideal case it can help predict whether neighborhoods having a certain set of attributes (most popular venues in this case), will draw a certain type of crimes more than others.

So, let's see and dig into the problem and figure out what do we need to start this project.

Data Sources:

1. Philadelphia Police Department crimes data set. It includes both violent and nonviolent offences. It has location, general crime type, datetime and latitude and longitude of that address <https://www.opendataphilly.org/dataset/crime-incidents>
2. Foursquare API – similar to the course material, I will use it information about neighborhoods
3. Neighborhood shape data in GeoJson format. This will be used to display neighborhoods borders on the map and also latitude and longitude of crime will be compared vs shapes from this source to map each crime into appropriate neighborhood. https://github.com/azavea/geo-data/tree/master/Neighborhoods_Philadelphia

Methodology:

I used the datasets from above to generate datasets ready to be used for clustering analysis.

DATA CLEANSING:

- Neighborhood Crime dataset: Just a regular cleansing with datatypes/ categorical types. The most cumbersome task was to find a proper neighborhood given latitude and longitude of a crime. I used function borrowed from the link in the code to run through all 3M records of crime data and assign each record to a neighborhood based on these coordinates. It was a very long process and it took more than 8 hours to process all rows on a standard “cognitive labs” python virtual server.
Next, I encoded crime type as a one hot encode and processed to that data set to standard scikit learn k-means clustering algorithm.
- Neighborhood venues dataset. Also, just a regular cleansing similar to the course workbook. The original dataset was obtained through Foursquare API as required by the course rubric.
Next, one hot encoded data was feed into clustering algorithm.

CLUSTERING

I used K-means clustering method for analyzing for patterns and grouping. To come up with the optimal clusters amount I used Elbow approach for both crime rates and neighborhood analysis. The inertia amount was starting to show less degree of declining with cluster amounts of 10 for both subsets. Hence I choose 10 as an optimal N for both datasets.

Results:

The final dataset consists of each neighborhood clustered by venues and crimes. Also, each line includes top ten venues and top ten crimes. During EDA, few neighborhood clusters showed clear relationship with definite crime clusters, while majority of them had at least 2-3 crime clusters without and dominant ones.

For example, neighborhoods clusters under 0,3 and 6 cluster only have crimes clustered 0,7 and 8 respectively.

Some neighborhoods (cluster 2,4,5) showed predominant crime type (crime clusters 5,6), while roughly half of neighborhood clusters had no clear dominant clustering in crime.

Let's try to determine if there are any similarities in clusters and whether we can spot any trends. Again, we will focus only on Neighborhood clusters that only have 1 or 2 predominant crime types (0,3,6) in it, ignoring neighborhoods with all sort of crimes.

Crime cluster 0 is a cluster with predominant all other offences and thefts from vehicles. Then Thefts, Vandalism/Fraud are trailing behind. While all of this is not really great, still it is much more "safer" than Cluster crime 7 for example, where "All other assaults" are taking 1 and the second most frequently happening place.

Let's see if we draw any conclusions: Neighborhood cluster 0 which draws so many "light" crimes is a cluster where many service establishments are located: Home Service/Vintage Store/Barbershop/Zoo Exhibit – all nonfood places with some exceptions. Accordingly, these spots with mass gatherings, though not as social as eating/drinking draw many "light crimes". On the contrary, Neighborhood cluster 6 is more of "eatery" style borough with Thai/Fast Food/Filipino restaurants. It looks like that these (not necessary these specific types of restaurants, but more like general places where people can get food and drinks) – tend to attract more violent type of crimes.

Discussion:

Clearly there is not enough information gained by using this data and trying to see how clusters obtained for both crime and venues dataset relate to each other. We can observe that small bars and eateries can relate to more serious crimes vs neighborhoods that have more like public gathering venues, where predominantly other offences/thefts occur mostly. Still we identified only few clusters out of 10 where this relationship has a clear effect. For the rest – we can't say with a great certainty what type of crime will prevail in that specific neighborhood and this is clearly something to work on in future developments.

Also, there can be many other steps to improve model accuracy (specifically on the crimes side, as there are too many "Other" types of crimes).

Additionally, some tweaking with venues data is also possible: for example, aggregate some venues into major categories (ex. Small restaurants, bars, etc.). This can improve understanding of clusters since now there a lot of features in one hot encoded field and it makes interpreting results somewhat complicated.

Conclusion:

In the final project I applied the approach discussed during the class to analyze city neighborhoods. I decided to work on city of Philadelphia neighborhoods. Additionally, as I was interested in seeing whether there is a relationship between clusters of venues types and crime types, I applied similar K-means clustering technique on Philadelphia police dataset and came up with neighborhoods split by crime cluster.

In the end out of 10 clusters I was able to relate 3 definite venue clusters with specific crime rates clusters (see results section). The other 7 clusters had some mixed results without any predominant crime cluster.