

Team 2

Chicago Car Accidents Assessment Utilizing Clustering Analysis

1st Hodgetts, Michael
Electrical Engineering Dept
EE-695 - Machine Learning
Hoboken, NJ
mhodgetts@comcast.net

2nd Paladugu, Rithvika
Computer Engineering
EE 695 - Machine Learning
Hoboken, NJ
paldugurithvika@gmail.com

Abstract—Practice as a team to analyse unsupervised data by exploring various clustering methods (k-means, k-modes, kprototype, Gaussian and DBSCAN) on Chicago Accident database.

I. INTRODUCTION

Utilizing two datasets from Chicago Department website (Crashes and People related to Crashes), we like to find areas/neighborhoods within the city that have different characteristics in terms of the attributes available. Dataset is large (750K rows) so we decided to look at only 2 years, i.e. 2021 and 2022 which reduces the row count to about 250K. There is a good amount of prep to get data suitable for cluster analysis (ETL, Hot Encoding, data cleanup) We tried 5 different cluster method with varying results. We will illustrate some of the main differences these assigned clusters present. Provide a summary of each cluster characteristics to illustrate the main differences between them.

II. RELATED WORK

There are many examples on how clustering solves problems in various business use cases. Segmentation analysis, anomaly detection and a form of classification are some of the use cases. There are also Deep Learning techniques that we could have explored but we ran out of time. The analysis here could be useful to identify hot areas of accidents due to poor signage, road conditions, and other conditions that could be improved to reduce accidents or injuries in the Chicago area. Practicing how to get meaningful insights is essential in the business world.

III. OUR SOLUTION

Since the data is large and has many attributes (both numerical and categorical), reducing the number of attributes will be critical. Also, minimizing the number of clusters needed is beneficial to ensure each cluster has meaningful differences and similar volume. Heatmaps showing differences will be provided to help show major allocation of attributes to each cluster. Another issue is we have a lot of category data as well as continuous. k-modes was explored first since

it handles categorical data and kmeans covers numerical. K-Prototypes covers both at the same time. We tried to use Gaussian Mixture but that seems to not be appropriate and gave us unusual results. Another method is DBSCAN which can handle hot encoding but so far find it difficult to gain the necessary clusters needed, very sensitive to (EPS) and size of data. kprototype so far has the best results with the mixed data. We will show Heatmaps to illustrate our results and provide summary description of the final clustering results. Significant code related to grouping data to assess the cluster allocation was necessary in evaluating key differences for each cluster.

A. Description of Dataset

The dataset can be found here: Chicago Crashe Data, Two datasets; one for crashes in Chicago and the other are the people characteristics of those crashes. The 'Crash ID' is the key attribute to join the 2 datasets. We reduced the size to only years 2021 and 2022 which still leaves about 250K rows and 68 columns. We used Knime (ETL software) to help do this join and ensure it was properly achieved. The output of the Knime workflow is our starting point for the python code section. Here is the workflow using KNIME ETL software:

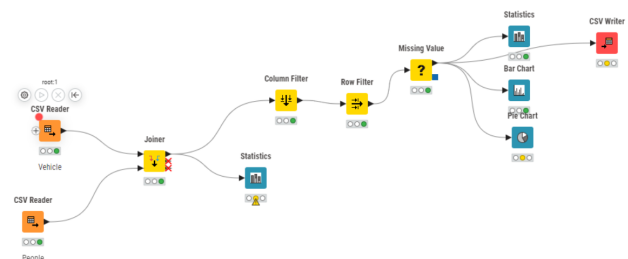


Fig. 1. Chicago Dataset Prep

Data had very few problems with missing data but we removed as necessary or added average values as needed. We separated the dataset into two parts, categorical and numerical. Performed Hot Encoding on the categorical and scaling on the numerical (0,1), then combined them back together as one

file. New Shape of the file is now 248K by 111 columns with hot encoding. We created another dataset for k-prototypes which has 31 columns based on what attributes we deemed valuable. Pie charts were useful in this particular dataset since most attributes are categorical. Many Attributes had too many categories that added more complexity so therefore we re-classify those less than 1 percent into one group called remainder to help reduce the number of hot columns needed.

Below is the list of columns we are working with in the dataset. Some have been removed since they will not contribute to the model performance.

```
df.columns
Index(['CRASH_RECORD_ID', 'RD_NO', 'CRASH_DATE', 'POSTED_SPEED_LIMIT',
      'TRAFFIC_CONTROL_DEVICE', 'DEVICE_CONDITION', 'WEATHER_CONDITION',
      'LIGHTING_CONDITION', 'FIRST_CRASH_TYPE', 'TRAFFICWAY_TYPE',
      'ALIGNMENT', 'ROADWAY_SURFACE_COND', 'ROAD_DEFECT', 'REPORT_TYPE',
      'CRASH_TYPE', 'INTERSECTION_RELATED_I', 'NOT_RIGHT_OF_WAY_I',
      'HIT_AND_RUN_I', 'DAMAGE', 'DATE_POLICE_NOTIFIED',
      'PRIM_CONTRIBUTORY_CAUSE', 'SEC_CONTRIBUTORY_CAUSE', 'STREET_NO',
      'STREET_DIRECTION', 'STREET_NAME', 'BEAT_OF_OCCURRENCE', 'NUM_UNITS',
      'MOST_SEVERE_INJURY', 'INJURIES_TOTAL', 'INJURIES_FATAL',
      'INJURIES_INCAPACITATING', 'INJURIES_NON_INCAPACITATING',
      'INJURIES_REPORTED_NOT_EVIDENT', 'INJURIES_NO_INDICATION',
      'INJURIES_UNKNOWN', 'CRASH_HOUR', 'CRASH_DAY_OF_WEEK', 'CRASH_MONTH',
      'LATITUDE', 'LONGITUDE', 'LOCATION', 'PERSON_ID', 'PERSON_TYPE',
      'CRASH_RECORD_ID (right)', 'RD_NO (right)', 'VEHICLE_ID',
      'CRASH_DATE (right)', 'CITY', 'STATE', 'ZIPCODE', 'SEX', 'AGE',
      'DRIVERS_LICENSE_STATE', 'DRIVERS_LICENSE_CLASS', 'SAFETY_EQUIPMENT',
      'AIRBAG_DEPLOYED', 'EJECTION', 'INJURY_CLASSIFICATION', 'HOSPITAL',
      'EMS_AGENCY', 'DRIVER_ACTION', 'DRIVER_VISION', 'PHYSICAL_CONDITION',
      'PEDPEDAL_ACTION', 'PEDPEDAL_VISIBILITY', 'PEDPEDAL_LOCATION',
      'BAC_RESULT', 'BAC_RESULT VALUE'],
      dtype='object')
```

Fig. 2. Chicago Dataset

EDA

Here we show a sampling of some of the attributes classes for the categorical data and numerical. There are many columns so we will just show a sample of some. You can see more in the code readout.

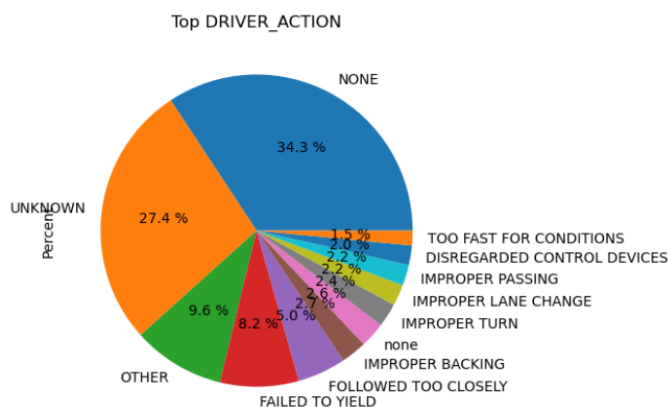


Fig. 3. Categorical Pie Charts (TOP Driver Reasons)

Here are some of the numeric attributes to examine note:

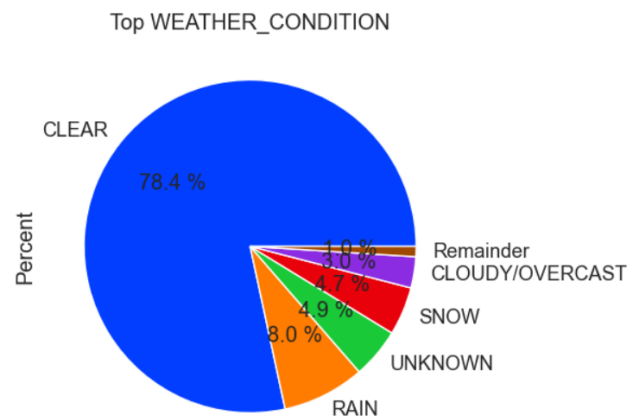


Fig. 4. Categorical Pie Charts (Sampling View)

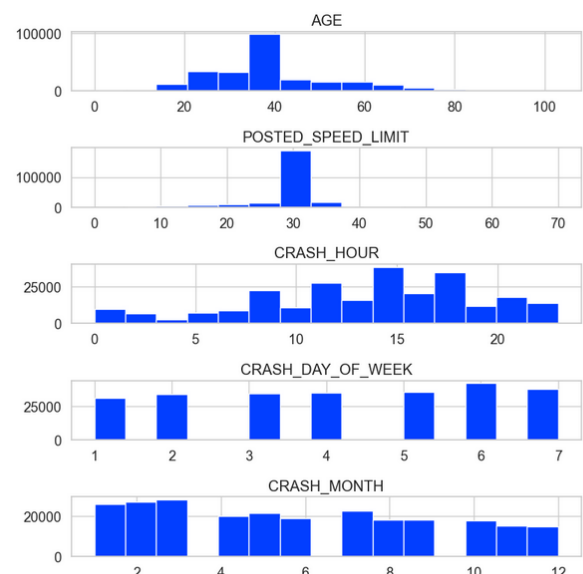
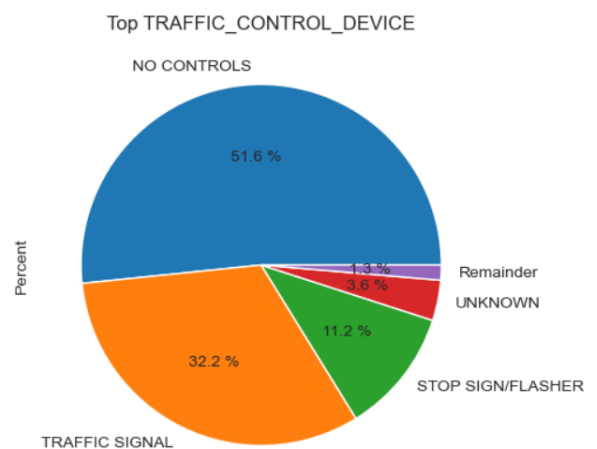


Fig. 6. Numeric Histograms (Sampling View)

B. Machine Learning Algorithms

Our data is raw and has no classification or specific purpose so it lends itself to utilize unsupervised data techniques. We explored K-modes, kmeans, kprototype and DBSCAN and will investigate other possible methods to find insights. kprototype is our best hope so far to get good results since it handles both categorical and numeric attributes. Kmodes can only handle category and kmeans does numerical. Noticed Kprototype takes a very long time to process the model. Elbow curve took over 24 hours. Since we have a large dataset, training take a good length of time for all model types. We created elbow curves for both K-modes and Kprototypes in the figures below:

From figures 2 and 3 you can see 4 clusters would be reasonable. We currently picked 4 clusters as a start. The difficulty is assessing the cluster characteristics with so many attributes. We did some Chi Testing on the attributes to find any that were not relevant but all picked were significant. We utilized pspark to use groupby by clusters by percent

	variable	chi2_test_stat	p_value	dof
34	INJURIES_UNKNOWN	0.000000e+00	1.000000e+00	0
41	PERSON_ID	1.243115e+06	4.985663e-01	1243110
67	BAC_RESULT VALUE	5.929951e+02	2.216402e-39	205
45	VEHICLE_ID	1.219355e+06	1.280826e-45	1197360
36	CRASH_DAY_OF_WEEK	8.620735e+02	5.757343e-162	30
...
30	INJURIES_INCAPACITATING	1.071638e+04	0.000000e+00	30
31	INJURIES_NON_INCAPACITATING	5.991988e+04	0.000000e+00	55
32	INJURIES_REPORTED_NOT_EVIDENT	2.847958e+04	0.000000e+00	40
33	INJURIES_NO_INDICATION	4.908476e+04	0.000000e+00	135
68	Cluster	1.243115e+06	0.000000e+00	25

Fig. 7. Chi Square Testing View

of occurrence to see what patterns emerge. We may have to build this view outside of Python to get a good illustrative view with excel. At the current time, we're assessing how to describe these clusters and find recommendations to provide back to the city of Chicago. This figure below illustrates the comparison we are doing with 4 clusters versus the categorical columns. You can see clearly differences in of how the clustering divided up the percent allocation. Percent allocation is based total population so we can compare to each different attribute. Will do the same for the numerical but it will look a little different since we will use averages.

C. Implementation Details

This subsection describes details of your implementation. Please focus on how you test and validate the performance, tune the hyperparameters, and select the best-performing models. Elaborate on techniques that you apply to improve the performance and explain why you use these techniques. You include few most important results/figures to illustrate your idea but do not let figures/tables dominate the content of the report. You can include few lines of critical code if needed.

But please avoid paste lengthy code in your report. Please make sure the figures/tables/code snapshots are of appropriate size including the font size.

IV. COMPARISON

In this section we'll discuss the results of each method and show more in detail the best solution results. The table below shows a summary of the different clustering algorithms.

Team 2 Model Summary					
Method	Application	Model Training Speed	Cluster Volumes	Error	Results/Notes
kmodes	Categorical	7.5 minutes	1) 32.1% 2) 30.4% 3) 23.4% 4) 14%	1375000	kmodes and kprototypes were similar in categorization but kmode had higher error
kmeans	Numerical	1.8 seconds	1) 25% 2) 25% 3) 25% 4) 25%	1 x10 ¹⁵	equal distribution but can only be used for numerical
GMM	Both	5 seconds	1) 58% 2) 28.2% 3) 11.8% 4) 1.7%	AIC=-1.33x10 ⁷	One cluster dominated in the distribution, results were not as good as KP
DBSCAN	Numerical	5.5 seconds	1) 70.3% 2) 17.5% 3) 6.3% 4) 6%		One cluster dominated in the distribution and results didn't look reasonable
Kprototype	Both	51 minutes	1) 28.4% 2) 27.4% 3) 23.5% 4) 20.6%	128000	Best of All, Reasonable results but took much longer to train

Fig. 8. Model Summary

Kmodes

Kmodes is designed for categorical data only and it can take the data without hot encoding. K-modes was successful and did a good job of clustering but you have to take out the numerical data and run that on kmeans. The cluster allocation is similar to K-prototype and its error was also one of the lowest in term of SSE.

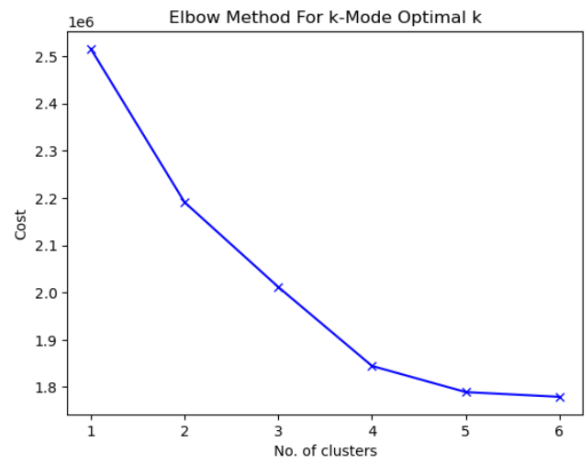


Fig. 9. k-mode

Kmeans

Kmeans above in the chart reflects the numerical data only, you can try to run the categorical data with hot encoding but its accuracy or ability to cluster diminishes.

Here is the numerical cluster view, We used mean and count to assess the numerical data

The numerical results were disappointing in that the shared volume across the clusters is similar for larger allocations but there are some minor differences between for the smaller

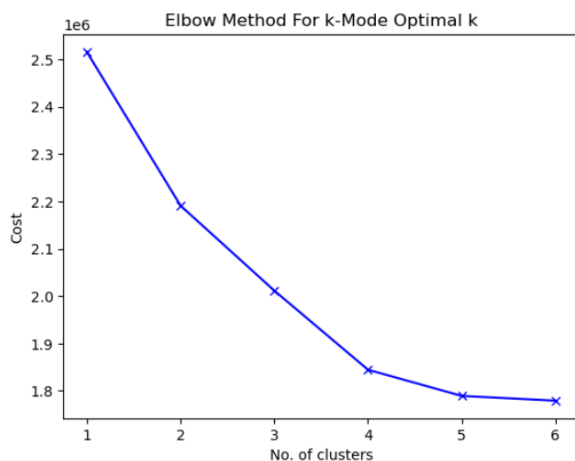


Fig. 10. k-mode

Sum of perc_of_count_total	Column	0	1	2	3
Row Labels					
ALIGNMENT		69.9264	6.234	5.924	17.3648
Remainder		0.706	0.0756	0.0628	0.162
STRAIGHT AND LEVEL		68.2392	6.086	5.7756	16.9188
STRAIGHT ON GRADE		0.9812	0.0724	0.0856	0.284
CRASH_TYPE		69.9264	6.234	5.924	17.3648
INJURY AND / OR TOW DUE TO CRASH		22.4772	2.044	1.9196	5.3072
NO INJURY / DRIVE AWAY		47.4492	4.19	4.0044	12.0576
DAMAGE		69.9264	6.234	5.924	17.3648
\$500 OR LESS		6.5468	0.5324	0.5108	1.5152
\$501 - \$1,500		16.9176	1.4428	1.398	4.0524
OVER \$1,500		46.462	4.2588	4.0152	11.7972
DEVICE_CONDITION		69.9264	6.234	5.924	17.3648
FUNCTIONING PROPERLY		27.3164	2.5424	2.4064	6.7736
NO CONTROLS		37.2568	3.1744	3.0032	9.0824
Remainder		1.0532	0.0994	0.0928	0.2524
UNKNOWN		4.3	0.4268	0.4216	1.2564
FIRST_CRASH_TYPE		69.9264	6.234	5.924	17.3648
ANGLE		9.3232	0.922	0.8264	2.3432
FIXED OBJECT		2.1652	0.2068	0.196	0.5216
HEAD ON		0.778	0.0688	0.0624	0.1896
PARKED MOTOR VEHICLE		10.5436	0.8124	0.7348	2.4116
PEDALCYCLIST		1.128	0.0748	0.0528	0.1656
PEDESTRIAN		1.6504	0.1652	0.1904	0.4004
REAR END		17.298	1.5396	1.4956	4.3836
REAR TO FRONT		1.182	0.108	0.104	0.3048
Remainder		1.5448	0.146	0.12	0.3908
SIDESWIPE OPPOSITE DIRECTION		1.1624	0.108	0.0856	0.2884
SIDESWIPE SAME DIRECTION		11.4648	1.0072	1.0308	2.898
TURNING		11.686	1.0752	1.0252	3.0672

Fig. 12. GMM Cluster Results

insignificant occurrences. Below we show an example for Posted Speed limit allocation for kmeans, you see very little distinction among the clusters.

POSTED_SPEED_LIMIT	62388	62110	62120	62005	248623
0	110	112	114	109	445
1			2		2
2	4			2	6
3	19	20	15	18	72
4				1	1
5	193	162	163	132	650
7		5			5
9	1		2	2	5
10	1438	1295	1394	1242	5369
11			1	3	4
12				1	1
14			1		1
15	2065	1710	1812	1699	7286
20	2425	2408	2354	2395	9582
22			1		1
24	6	2	5	11	24
25	3816	3594	3555	3806	14771
29		2			2
30	46747	47290	47105	47040	188182
32	2			2	4
33				2	2
34	6		4	2	12
35	4441	4072	4209	4341	17063
38			2		2
39	8	5	7	3	23
40	672	773	777	631	2853
45	354	542	530	479	1905
50	14	51	12	35	112
55	63	57	53	49	222

Fig. 11. kmeanexample for Posted Speed Limit

GMM

We used both numeric scaling and hot encoding here but we think that GMM is suited better with only numerical data. Here the cluster allocation was skewed towards one cluster that has the majority of volume. This was not our ideal solution. Here below is partial view of the cluster allocation. The first cluster is dominate in size with over 69 percent of the population.

DBSCAN

DBSCAN was very sensitive to the amount of data and columns. Only numerical data was used here and finding the

EPS and min sample values was a challenge, Had to iterate through 50 different EPS values to find a sweet spot. By changing EPS by only 0.01 increments at that sweet spot changed clusters size and number of cluster significantly. The best EPS value was 28

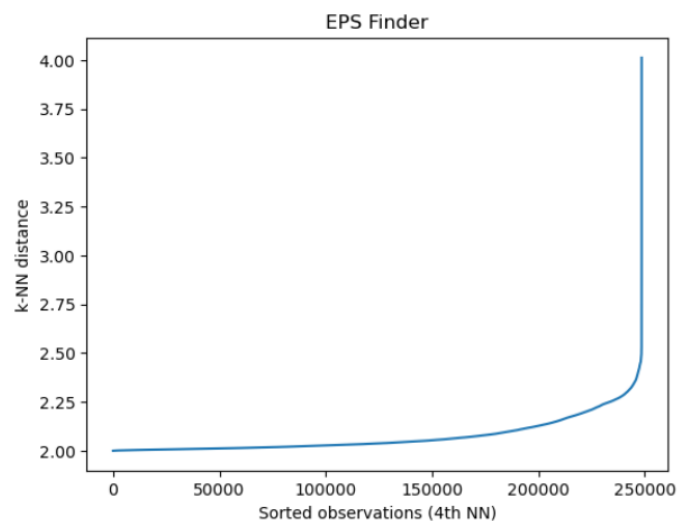


Fig. 13. DBSCAN EPS Finder

We ran category data as well and found it very sensitive to the size of data and number of columns, could not run with the required column size. Had to reduce the size of the file to get the results below (Show partial view):

The results were not optimal and we saw again high allocation to one cluster. The -1 or noise cluster had most of the volume, I spent a consider amount of time to get this result, very sensitive to parameter settings. We don't think DBSCAN was useful in this applicaiton.

K-prototype

K-Prototype was our best soluton in that we can use all the data and train at the same time. This took significantly

Sum of perc_of_count_total	Column Labels			
Row Labels	-1	0	1	2
ALIGNMENT	21.3076	18.4684	0.142	0.082
Remainder	0.3608	0.0396		
STRAIGHT AND LEVEL	20.4712	18.3384	0.142	0.082
STRAIGHT ON GRADE	0.4756	0.0904		
CRASH_TYPE	21.3076	18.4684	0.142	0.082
INJURY AND / OR TOW DUE TO CRASH	7.092	5.6008	0.002	8E-04
NO INJURY / DRIVE AWAY	14.2156	12.8676	0.141	0.081
DAMAGE	21.3076	18.4684	0.142	0.082
\$500 OR LESS	2.7968	0.9624	0.002	
\$501 - \$1,500	5.6228	4.114	0.024	0.004
OVER \$1,500	12.888	13.392	0.116	0.077
DEVICE_CONDITION	21.3076	18.4684	0.142	0.082
FUNCTIONING PROPERLY	9.1664	6.3196	4E-04	
NO CONTROLS	9.4284	12.0088	0.141	
Remainder	0.5308	0.0696		
UNKNOWN	2.182	0.0704	8E-04	0.082
FIRST_CRASH_TYPE	21.3076	18.4684	0.142	0.082
ANGLE	2.972	2.266		0.006
FIXED OBJECT	0.9584	0.3368		4E-04
HEAD ON	0.34	0.1392		
PARKED MOTOR VEHICLE	2.618	3.5136	0.138	0.004
PEDALCYLIST	0.4468	0.1092		
PEDESTRIAN	0.7644	0.1216		
REAR END	4.5964	5.2964		0.043
REAR TO FRONT	0.5656	0.1424	8E-04	0.002
Remainder	0.6656	0.2484		4E-04
SIDESWIPE OPPOSITE DIRECTION	0.4608	0.2452	0.001	8E-04
SIDESWIPE SAME DIRECTION	3.1868	3.2452	0.003	0.022

Fig. 14. DBSCAN Cluster Result

longer to train as you can see compared to others in the table. The volume allocation was similar to kmodes. So we can use kmeans for numerical and kmodes for categorical and combine them but KP does it together. We'll use KP to explain the main difference we see in the 4 clusters we trained on

linebreak

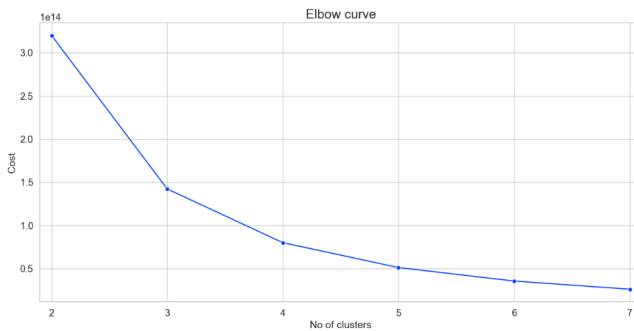


Fig. 15. k-prototype

Cluster-Characteristics

Below is a partial view of the KP Heatmap for the 4 clusters
Cluster 1 -In general 1 has the lowest occurrence of crash attributes -lowest occurrence of any damage level -Volume is lowest -unlikely for Hit and run -first 1-3 day of week occurrence

Cluster 2 -In general 2 has the second lowest occurrence of crash attributes -Likely to see Parked Vehicle Crash -likely Hit and Run -overindex on Sex X -overindex on Sex F -first 1-3 day of week occurrence

Cluster 3 -Overindex on Crash Injury -Overindex turning related crash type -Overindex on Intersection related crash -Overindex incapacitating injury -Overindex disregarding traffic

Sum of perc_of_count_total	Column Labels				Grand Total			Note
Row Labels	0	1	2	3				
ALIGNMENT	17.1056	26.1112	25.3148	30.9176	99.4492	min	max	major change
Remainder	0.7212	0.1844	0.2132	0.3968	1.0054	2	4	0.388354
STRAIGHT AND LEVEL	16.0024	25.6056	24.4892	30.1296	97.0108	1	4	0.232093
STRAIGHT ON GRADE	0.2932	0.318	0.4164	0.3956	1.4232	1	3	0.167097
CRASH_TYPE	17.1056	26.1112	25.3148	30.9176	99.4492	1	4	0.230585
INJURY AND / OR TOW DUE TO CRASH	5.0084	4.3912	16.3424	6.006	31.748	2	3	0.710969 xx
NO INJURY / DRIVE AWAY	12.0972	21.72	8.9724	24.9116	67.7012	3	4	0.449096
DAMAGE	17.1056	26.1112	25.3148	30.9176	99.4492	1	4	0.230585
\$500 OR LESS	3.48	3.4884	1.8864	2.8004	9.1052	1	4	0.21836
\$501 - \$1,500	4.25	7.0856	3.7328	8.7424	23.8108	1	4	0.398638
OVER \$1,500	10.9956	16.5672	19.6956	19.2748	66.5332	1	3	0.24084
DEVICE_CONDITION	17.1056	26.1112	25.3148	30.9176	99.4492	1	4	0.230585
FUNCTIONING PROPERLY	5.496	5.0384	21.6872	6.8172	39.0388	2	3	0.818401 xx
NO CONTROLS	10.3228	18.8628	1.8648	21.4664	52.5168	3	4	0.677187 xx
Remainder	0.2916	0.3168	0.4916	0.3888	1.4888	1	3	0.240796
UNKNOWN	0.9952	1.8932	1.2712	2.2452	6.4048	1	4	0.356245
FIRST_CRASH_TYPE	17.1056	26.1112	25.3148	30.9176	99.4492	1	4	0.230585
ANGLE	2.054	2.2632	6.3104	2.7872	13.4148	1	3	0.594899 xx
FIXED OBJECT	0.6588	0.666	0.5636	1.2012	3.0896	3	4	0.375001
HEAD ON	0.3988	0.266	0.316	0.318	1.0988	1	4	0.203957
PARKED MOTOR VEHICLE	1.5084	10.1572	0.3744	2.4624	14.5028	3	2	1.223868 xxx
PEDALCYLIST	0.2192	0.1828	0.6348	0.3844	1.4212	2	3	0.579652
PEDESTRIAN	0.4892	0.4324	1.1116	0.3732	2.4064	4	3	0.570616
REAR END	5.0976	3.3412	4.2936	11.9844	24.7168	2	4	0.637001
REAR TO FRONT	0.3344	0.4976	0.1596	0.7072	1.6988	3	4	0.549776
SIDESWIPE OPPOSITE DIRECTION	0.533	0.5256	0.244	0.894	2.2016	3	4	0.483789
SIDESWIPE SAME DIRECTION	0.344	0.5164	0.2248	0.5592	1.6444	3	4	0.377447
TURNING	3.1616	4.9012	2.3192	6.0188	16.4008	3	4	0.407513
HIT_AND_RUN_I	17.1056	26.1112	25.3148	30.9176	99.4492	1	4	0.230585
N	0.1872	0.4448	0.3756	0.3624	1.39	1	2	0.320322
none	14.7576	6.5256	19.2136	36.7872	67.284	2	4	0.503631
Y	2.1608	19.1408	5.7256	3.748	30.7752	1	2	1.009827 xx
INTERSECTION_RELATED_I	17.1056	26.1112	25.3148	30.9176	99.4492	1	4	0.230585
N	0.2192	0.2308	0.4376	0.3072	1.1948	1	3	0.336432
none	13.5644	23.7784	6.3476	28.1152	71.8056	3	4	0.548782
Y	3.322	2.102	18.5296	2.4952	26.4488	2	3	1.204015 xx

Fig. 16. Example of Cluster Heatmap for Categorical Attributes

signals -Overindex device condition functioning properly -
Overindex disregarding stop sign -Overindex likely to report on scene -likely to occur in the first 6 months of years

Cluster 4 -Volume is highest -Overindex Parked Motor Vehicle -Overindex on no injury -Overindex no control for device condition (non traffic light) -Overindex to be rear ended -Overindex for follow to closely crash type -Overindex for months 7-12

V. FUTURE DIRECTIONS

We still need to try other methods and see if the allocations are significantly different. We'll describe each cluster characteristics in more detail. We like to show so mapping features of how the clusters map over Chicago but find this difficult so far to accomplish. There another method call Squeezer which can be used with mixed data but has little documentaion. Deep learning techniques utilaing autoencoders can be examined. We'll add to this section for our final report. Ho to assess cluster info could be ahot topic

VI. CONCLUSIONS

We'll add our conclusions at the final report

REFERENCES

- 1) <https://scottmduda.medium.com/categorical-clustering-of-pittsburgh-car-accidents-using-k-modes-7c842cc15d87>
- 2) <https://data.cityofchicago.org/Transportation/Traffic-Crashes-Crashes/85ca-t3if>
- 3) <https://www.freecodecamp.org/news/8-clustering-algorithms-in-machine-learning-that-all-data-scientists-should-know/>
- 4) https://scikit-learn.org/stable/auto_examples/cluster/plot_abscan.html
- 5) <https://antonruberts.github.io/kproto-audience/>