

## Checkpoint 4: Machine Learning

### Introduction & Data preparation

In this checkpoint, we will use the repeaters' data from 2014-2018. For each repeater in this time period, we keep the count of each officer's misconduct under each misconduct category. Therefore a sample data will look like this:

	total	Conduct Unbecoming (Off-Duty)	False Arrest	Illegal Search	Lockup Procedures	Operation/Personnel Violations	Use Of Force	Verbal Abuse
60	25.0	1.0	2.0	10.0	2.0	2.0	8.0	0.0
72	3.0	0.0	0.0	1.0	0.0	0.0	1.0	1.0
135	11.0	0.0	3.0	2.0	1.0	2.0	3.0	0.0
152	31.0	0.0	0.0	1.0	1.0	13.0	9.0	3.0
193	3.0	0.0	0.0	0.0	0.0	1.0	1.0	1.0

The first column represents the officer ids, where subsequent columns give the total as well as subtotal of the 19 categories. We then convert all the subtotals to ratios between subtotal and total. The point is to make sure for each officer, now the subtotals will sum to 1 and we thus normalize the data as machine learning data input. The transformed data looks like this:

	total	Conduct Unbecoming (Off-Duty)	False Arrest	Illegal Search	Lockup Procedures	Operation/Personnel Violations	Use Of Force	Verbal Abuse
60	25.0	0.04	0.08	0.40	0.08	0.08	0.32	0.00
72	3.0	0.00	0.00	0.33	0.00	0.00	0.33	0.33
135	11.0	0.00	0.27	0.18	0.09	0.18	0.27	0.00
152	31.0	0.00	0.00	0.03	0.03	0.42	0.29	0.10
193	3.0	0.00	0.00	0.00	0.00	0.33	0.33	0.33

### Methods

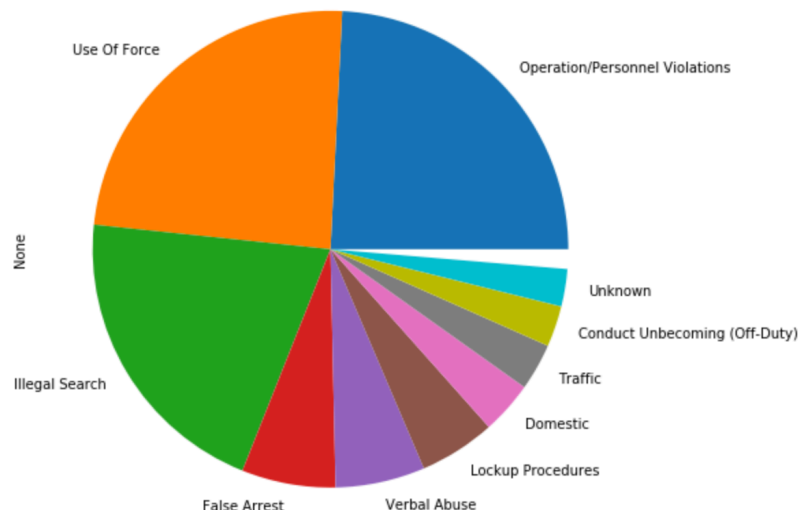
We use K-Means Clustering as our machine learning model. We first use s-score to determine the optimal number of clusters and we get optimal K equals 2. Then we train the model and generate label (0 or 1) for each officer record and then split the dataframe according to the generated label.

### Results

Cluster 0 consists of 79 officers out of a total of 557 officers (14% of the population). The mean count of total allegations is 60.65 with a standard deviation of 20. Therefore the 95% confidence interval constructed is [56.235, 65.056]

A more detailed breakdown for the ratio of each category v.s all categories is as follows:

Conduct Unbecoming (Off-Duty)	0.027848
False Arrest	0.063797
Illegal Search	0.206076
Lockup Procedures	0.051646
Operation/Personnel Violations	0.242152
Use Of Force	0.241392
Verbal Abuse	0.060886
Criminal Misconduct	0.020380
Domestic	0.035823
Traffic	0.031772
Bribery / Official Corruption	0.002911
Supervisory Responsibilities	0.008481
Money / Property	0.004557
Medical	0.000253
Racial Profiling	0.000253
Excessive Force	0.000000
Drug / Alcohol Abuse	0.000759



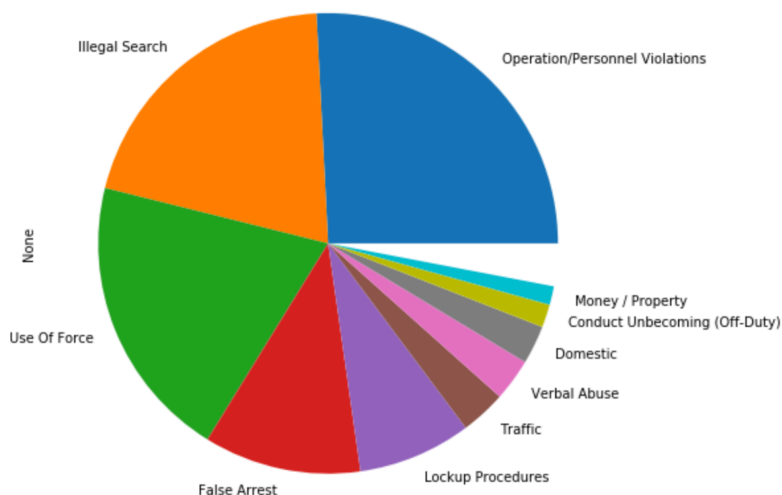
We can see that the major misconduct categories are:

- 1) Operations/Personnel Violations
- 2) Use of Force
- 3) Illegal Search
- 4) False Arrest
- 5) Verbal Abuse

Cluster 1 consists of 478 officers out of a total of 557 officers (86% of the population). The mean count of total allegations is 16.56 with a standard deviation of 9.40. Therefore the 95% confidence interval constructed is [15.718, 17.403]

A more detailed breakdown for the ratio of each category v.s all categories is as follows:

Conduct Unbecoming (Off-Duty)	0.016234
False Arrest	0.110439
Illegal Search	0.203138
Lockup Procedures	0.079728
Operation/Personnel Violations	0.258138
Use Of Force	0.200795
Verbal Abuse	0.029791
Criminal Misconduct	0.011130
Domestic	0.026987
Traffic	0.031841
Bribery / Official Corruption	0.005230
Supervisory Responsibilities	0.006820
Money / Property	0.013013
Medical	0.000377
Racial Profiling	0.001381
Excessive Force	0.001025
Drug / Alcohol Abuse	0.002762
Unknown	0.008368



We can see that the major misconduct categories are:

- 1) Operations/Personnel Violations
- 2) Illegal Search
- 3) Use of Force
- 4) False Arrest
- 5) Lockup Procedures

## Answers

From the confidence interval we calculated, we can project that the total count of these repeaters will be in the range of [11956, 13458]. Combining the graph with the ratio, We can find that officers will be more probable towards Operations/Personnel Violations, Illegal Search, Use of Force, False Arrest, Lockup Procedures and Verbal Abuse.

This Corroborates with our previous findings and we add Lockup Procedures and Verbal Abuse to our attention. The misconduct pattern, in terms of the categories, is pretty clear. Over 80% of the misconducts are accounted for by the above categories. I would say that Use of Force, False Arrest, Lockup are all misconducts that highly depend on the discretion of the officers. The varying of severity of scenarios, surroundings, officer-victim relationships, etc all complicate the analysis. I would say to look deeper into the issue, it would require the police department to be transparent and objective for the reports. The censor of reports prior to releasing can be something that hugely impacts our perspective regarding these issues.