

Dallas Police Public Data - RMS Incidents

Predicting the whether Drugs are involved in Crime

Sebastian, Jules	040 983 119
Simpson, Skyler	040 951 590

Data Collection

Source: Dallas Police Public Data - RMS Incidents

<https://www.dallasopendata.com/Public-Safety/Police-Incidents/qv6i-rri7>

How the data was collected

The data was sourced from preliminary reports given to the police after June 1st 2014.

Preprocessing

The Initial Data was Enormous, approx 875 000 Rows of Data, denoted by 86 Attributes.

- Lots of data does not mean the data is good.

Duplicate data was Removed

Redundant Attributes were Merged

Outliers were Removed

Correlating Attributes were filled based on the dataset

Missing Data led to removed rows: Where 6 or more attributes had missing values, the rows are ignored

Analysis

Classification

Vector distances are calculated between neighboring data points and classifying the dataset based on those distances.

kNearestNeighbors

We attempted kNN but a vector based approach was very costly, when nominal fields have such a large breadth of options.

- Long Calculations [>10 hours] and unreliable results

Analysis

Clustering

Grouping objects in such a way, that objects falling in a group would be more similar to those in other groups

- **Simple kMeans**
- **Farthest First**

Results

Clustering predicted approx 8% of incidents were drug related

With 10.85% of the data incorrectly clustered.

By projecting from the original dataset, we would expect the results to land closer to 3.85% of incidents are drug related.

This means our model overestimates the number of incidents.

kNN

1.91%

Time taken to test model on training data: 51823.52 seconds

=== Summary ===

Correctly Classified Instances	679686	98.1469 %
Incorrectly Classified Instances	12833	1.8531 %
Kappa statistic	0.6609	
Mean absolute error	0.0192	
Root mean squared error	0.098	
Relative absolute error	27.632 %	
Root relative squared error	52.6366 %	
Total Number of Instances	692519	

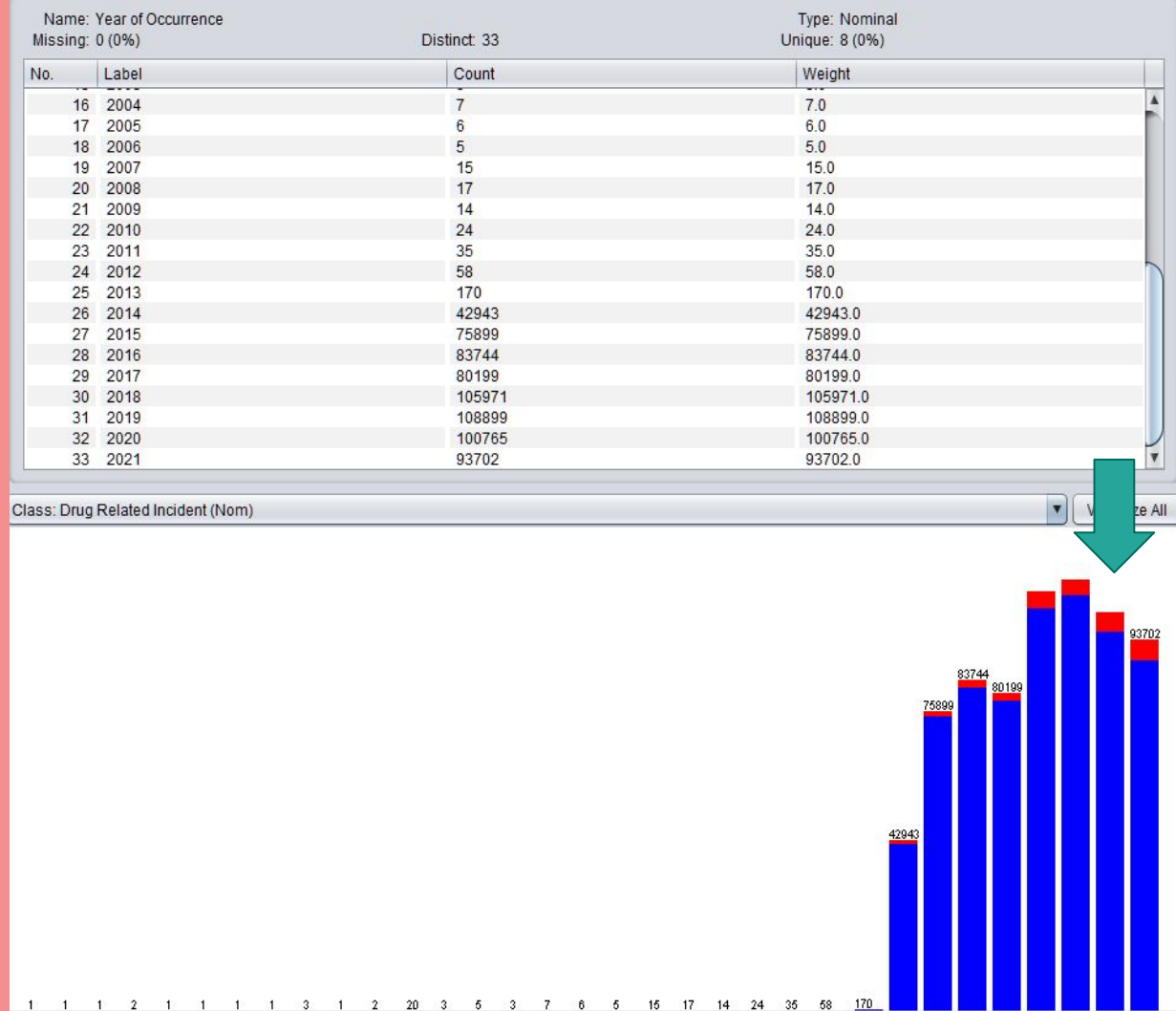
=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	MCC	ROC Area	PRC Area	Class
	0.999	0.477	0.982	0.999	0.990	0.690	0.995	1.000	No
	0.523	0.001	0.932	0.523	0.670	0.690	0.995	0.838	Yes
Weighted Avg.	0.981	0.460	0.981	0.981	0.979	0.690	0.995	0.994	

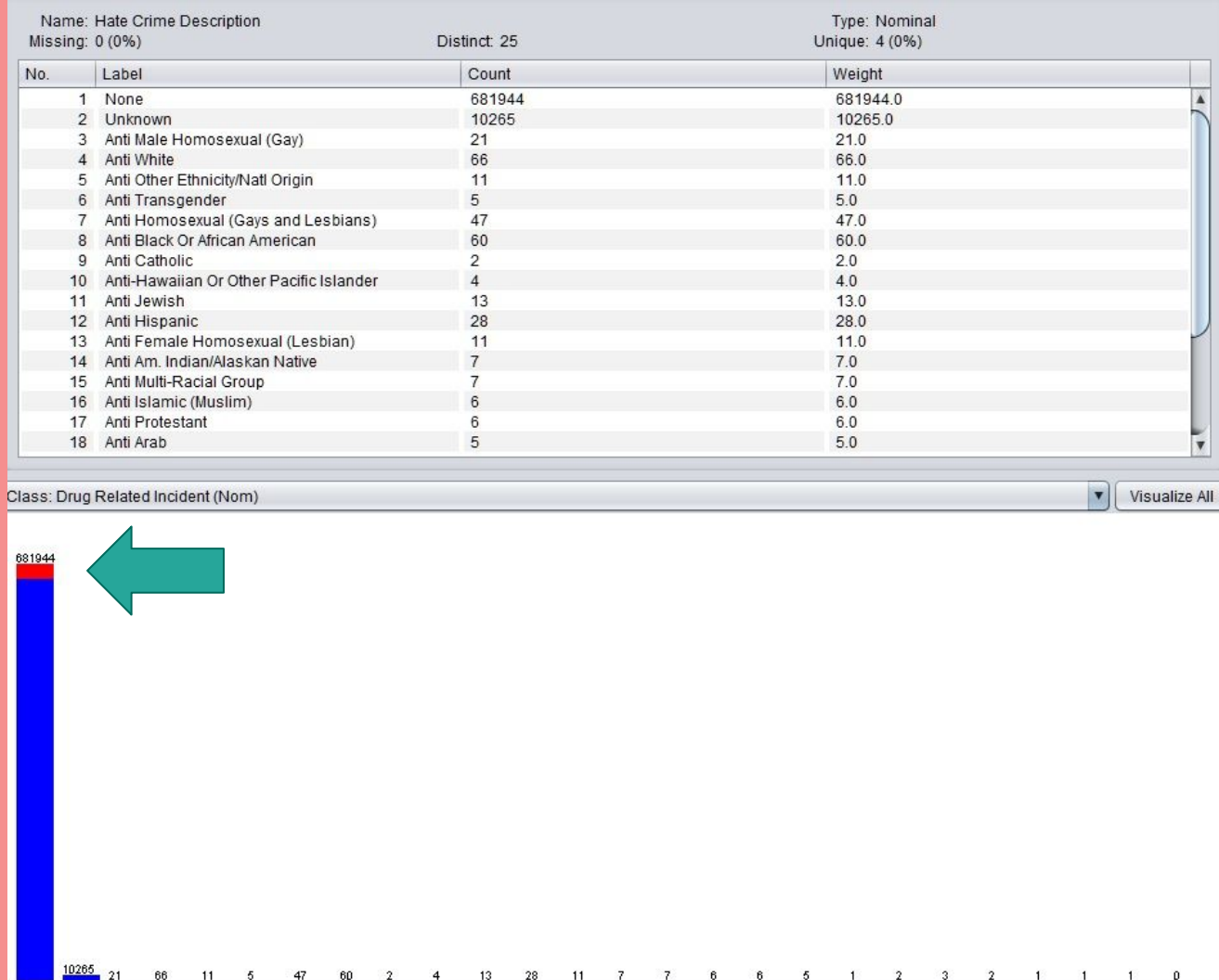
=== Confusion Matrix ===

a	b	<-- classified as
666678	946	a = No
11887	13008	b = Yes

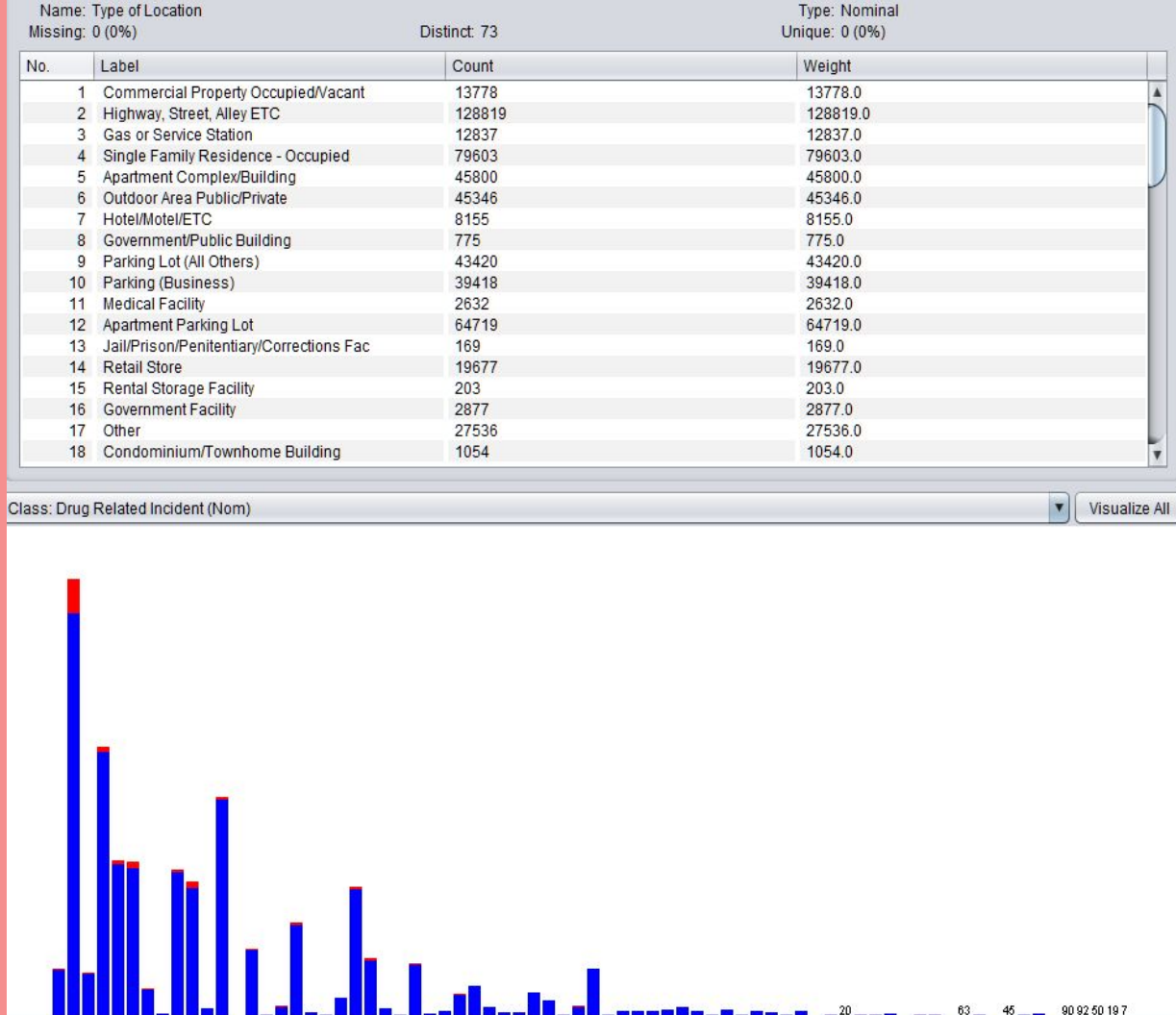
Drug Related Incidents are Increasing Over Time



Drugs and Hate Crimes Don't Mix



Most Drug
Related
Incidents
Occur On
Highways,
Streets or in
Alleyways



Conclusion

Our results had higher accuracy than we thought.

The data should be better cleaned.

The dataset contains too many sources of error.

Most of the information has no impact of our search

The data being populated is too unreliable for critical fields