

Binary classification:

* TO predict the severity of the crime binary classification is used.

* The required features for normalization is latitude, longitude, closest-station.

* Other features for the dataset are all categorical data & have all been converted into dummy variable (District, tgm-bldc, weekday, location description).

* We performed split on the 80,000 sample records, into training set & test set and normalized by themselves.

* Latitude & longitude info more helpful for classification, If the latitude & longitude is low, it's more likely that severe crime will happen.

* It is consistent with the fact that the safety of ^{ca}at south chigo is notoriously bad.

* The economic, unemployment rate, age status do provide the expected results. If for regions with lower income & higher unemployment rate, the crimes are going to be more severe.

* As a result if the tym blk is 3am to 6am, the proportion of severe crime is obviously higher, which is not ^{the} case in tym blk 9am to 12pm.

* The crime takes place in an appa apartment or house, it's more likely to be a severe crime. The crime takes place on a street where everyone can see, it's less likely going to be severe.

* Severe crime accounts for 46% of total crime and non-severe is around 53%.

* Built different classifiers using the training set, & examined the accuracy of all the classifiers using the test set.

* Based on the accuracy comparison of diff models, baseline model have least accuracy of 54%.

* If classification accuracy is the only criteria to judge whether a model is good or not, a good model have greater than the baseline model.

* I have compared all of the model logistic regression have the accuracy of 61%.

Severe

Arson, Assault, battery, ~~crim sexual assault~~
criminal damage, criminal trespass, homicide, robbery.