# A Final Project Report on Analysis of 'Chicago Crime' Dataset

**Bhavik Shah**

**Northeastern University**

**MS-IS (NU Id: 001825444)**

# Summary

This project is based on the Analysis of the 'Chicago Crime' dataset. The raw dataset consists of crime incidents taking place over Chicago. The data span a period of more than 17 years, including all ~100,0000 crime incidents from 2001 up to 2018. Crime Incident includes crimes category, Location of category, FBI code, cases resolved, and a plaintext review.

In the project, first the analysis of the raw dataset using Power BI is done. The **MapReduce** analysis is performed using several MapReduce patterns such as Partitioning, Binning, Job Chaining etc. Also, in this project, **Apache pig, Apache Hive and Apache Mahout** is implemented.

The analysis on the dataset can help to determine and predict the pattern of crime happening over the years and awareness and improvement of public safety.

Dataset Link: https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-present/ijzp-q8t2

## Chicago Crime Dataset:

| ID | Case Numl | Date | Block | IUCR | Primary Ty | Descriptio | Location D | Arrest | Domestic | Beat | District | Ward | Communit | FBI Code | X Coordina | Y Coordina | Year | Updated On | Latitude | Longitude | Location |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 8812461 | HV485448 | 9/6/2012 16:00 | 048XX N SHERIDAN RD | 842 | THEFT | AGG: FINA | RESIDENCI | FALSE | FALSE | 2024 | 20 | 48 | 3 | 6 | 1168762 | 1932302 | 2012 | 2/10/2018 15:50 | 41.96976 | -87.6548 | (41.969762615, -87.654835892) |
| 8812465 | HV485482 | 9/19/2012 12:30 | 088XX S CLYDE AVE | 1320 | CRIMINAL | TO VEHICL | STREET | FALSE | FALSE | 412 | 4 | 8 | 48 | 14 | 1191710 | 1846852 | 2012 | 2/10/2018 15:50 | 41.73476 | -87.5732 | (41.734755, -87.573238269) |
| 8812466 | HV485444 | 9/19/2012 12:30 | 005XX E 51ST ST | 890 | THEFT | FROM BUI | COMMER( | FALSE | FALSE | 223 | 2 | 3 | 38 | 6 | 1180467 | 1871309 | 2012 | 2/10/2018 15:50 | 41.80213 | -87.6137 | (41.802132606, -87.613677631) |
| 8812467 | HV485446 | 8/24/2012 9:00 | 001XX N CLARK ST | 890 | THEFT | FROM BUI | SCHOOL/ F | FALSE | FALSE | 111 | 1 | 42 | 32 | 6 | 1175522 | 1901045 | 2012 | 2/10/2018 15:50 | 41.88384 | -87.6309 | (41.883842668, -87.63092094) |
| 8812469 | HV485229 | 9/20/2012 20:00 | 106XX S MAY ST | 820 | THEFT | $500 AND | DRIVEWAY | FALSE | FALSE | 2232 | 22 | 34 | 73 | 6 | 1170655 | 1834224 | 2012 | 2/10/2018 15:50 | 41.70059 | -87.6507 | (41.700586268, -87.650741113) |
| 8812472 | HV485436 | 9/20/2012 15:15 | 051XX W 51ST ST | 1570 | SEX OFFEN | PUBLIC IN | CTA TRAIN | FALSE | FALSE | 814 | 8 | 23 | 56 | 17 | 1143102 | 1870328 | 2012 | 2/10/2018 15:50 | 41.80022 | -87.7507 | (41.800217334, -87.750736204) |
| 8812473 | HV485170 | 9/21/2012 9:14 | 099XX S WALLACE ST | 820 | THEFT | $500 AND | STREET | FALSE | FALSE | 2232 | 22 | 9 | 73 | 6 | 1174054 | 1838968 | 2012 | 2/10/2018 15:50 | 41.71353 | -87.6382 | (41.713529895, -87.638155035) |
| 8812475 | HV485454 | 9/21/2012 12:15 | 025XX W ALTGELD ST | 820 | THEFT | $500 AND | STREET | FALSE | FALSE | 1431 | 14 | 35 | 22 | 6 | 1159066 | 1916530 | 2012 | 2/10/2018 15:50 | 41.92669 | -87.6909 | (41.926688472, -87.690922996) |
| 8812476 | HV485458 | 9/21/2012 12:05 | 029XX S STATE ST | 1811 | NARCOTIC | POSS: CAN | CHA PARK | TRUE | FALSE | 133 | 1 | 3 | 35 | 18 | 1176760 | 1885391 | 2012 | 2/10/2018 15:50 | 41.84086 | -87.6268 | (41.840859176, -87.626847966) |
| 8812478 | HV485241 | 9/21/2012 9:30 | 013XX W 89TH ST | 486 | BATTERY | DOMESTIC | RESIDENCI | FALSE | TRUE | 2222 | 22 | 21 | 71 | 08B | 1169046 | 1845789 | 2012 | 2/10/2018 15:50 | 41.73236 | -87.6563 | (41.732357234, -87.656299678) |
| 8812479 | HV485398 | 2/1/2012 8:00 | 087XX S THROOP ST | 840 | THEFT | FINANCIAL | RESIDENCI | FALSE | FALSE | 2222 | 22 | 21 | 71 | 6 | 1169189 | 1846772 | 2012 | 2/10/2018 15:50 | 41.73505 | -87.6557 | (41.73505164, -87.655747456) |
| 8812484 | HV485120 | 9/21/2012 1:00 | 022XX W 111TH ST | 486 | BATTERY | DOMESTIC | RESIDENCI | FALSE | TRUE | 2212 | 22 | 19 | 75 | 08B | 1163198 | 1830971 | 2012 | 2/10/2018 15:50 | 41.69182 | -87.6781 | (41.691818411, -87.678136584) |
| 8812485 | HV485478 | 9/20/2012 18:00 | 060XX S KOMENSKY AVE | 1320 | CRIMINAL | TO VEHICL | VEHICLE N | FALSE | FALSE | 813 | 8 | 13 | 65 | 14 | 1150404 | 1864258 | 2012 | 2/10/2018 15:50 | 41.78342 | -87.7241 | (41.783421395, -87.724114952) |
| 8812486 | HV485412 | 9/21/2012 12:07 | 081XX S HONORE ST | 486 | BATTERY | DOMESTIC | RESIDENCI | FALSE | TRUE | 614 | 6 | 18 | 71 | 08B | 1165432 | 1850785 | 2012 | 2/10/2018 15:50 | 41.74614 | -87.6694 | (41.746144249, -87.669398156) |
| 8812487 | HV485427 | 9/21/2012 10:25 | 014XX S HARDING AVE | 2026 | NARCOTIC | POSS: PCP | STREET | TRUE | FALSE | 1011 | 10 | 24 | 29 | 18 | 1150280 | 1892777 | 2012 | 2/10/2018 15:50 | 41.86168 | -87.7238 | (41.861683877, -87.723827836) |
| 8812488 | HV484975 | 9/21/2012 3:35 | 002XX W ERIE ST | 031B | ROBBERY | ARMED: O | STREET | FALSE | FALSE | 1831 | 18 | 42 | 8 | 3 | 1174460 | 1904731 | 2012 | 2/10/2018 15:50 | 41.89398 | -87.6347 | (41.89398105, -87.034710466) |
| 8812489 | HV485466 | 9/21/2012 9:45 | 014XX S AVERS AVE | 2024 | NARCOTIC | POSS: HER | SIDEWALK | TRUE | FALSE | 1011 | 10 | 24 | 29 | 18 | 1150947 | 1892739 | 2012 | 2/10/2018 15:50 | 41.86157 | -87.7214 | (41.861566581, -87.72138037) |
| 8812490 | HV485408 | 9/21/2012 10:15 | 027XX N NORDICA AVE | 610 | BURGLAR\ | FORCIBLE | RESIDENC | FALSE | FALSE | 2512 | 25 | 36 | 18 | 5 | 1128618 | 1917505 | 2012 | 2/10/2018 15:50 | 41.92994 | -87.8028 | (41.92993546, -87.802785407) |
| 8812491 | HV485490 | 9/21/2012 9:00 | 101XX S ABERDEEN ST | 2825 | OTHER OF | HARASSM | RESIDENC | FALSE | TRUE | 2232 | 22 | 34 | 73 | 26 | 1170776 | 1837517 | 2012 | 2/10/2018 15:50 | 41.70962 | -87.6502 | (41.709620136, -87.650202406) |
| 8812495 | HV485402 | 9/21/2012 11:50 | 027XX N KILPATRICK AVE | 2093 | NARCOTIC | FOUND SU | SIDEWALK | TRUE | FALSE | 2521 | 25 | 31 | 19 | 26 | 1144588 | 1917701 | 2012 | 2/10/2018 15:50 | 41.93019 | -87.7441 | (41.930187136, -87.74409405) |
| 1964884 | HH146026 | 1/24/2002 17:59 | 010XX W 18 ST | 2022 | NARCOTIC | POSS: COC | SMALL RE] | TRUE | FALSE | 1233 | 12 | | | 18 | 1169515 | 1891579 | 2002 | 2/28/2018 15:56 | 41.858 | -87.6533 | (41.858000042, -87.65325445) |
| 8812496 | HV485420 | 9/20/2012 16:00 | 079XX S PEORIA ST | 320 | ROBBERY | STRONGAI | SIDEWALK | FALSE | FALSE | 621 | 6 | 17 | 71 | 3 | 1171682 | 1852294 | 2012 | 2/10/2018 15:50 | 41.75015 | -87.6465 | (41.750150526, -87.646452666) |
| 8812497 | HV485293 | 9/20/2012 22:00 | 048XX W LAKE ST | 870 | THEFT | POCKET-PI | CTA TRAIN | FALSE | FALSE | 1532 | 15 | 28 | 25 | 6 | 1144147 | 1901839 | 2012 | 2/10/2018 15:50 | 41.88667 | -87.7461 | (41.886668404, -87.746113681) |
| 8812498 | HV485439 | 9/21/2012 12:10 | 015XX N LOCKWOOD AVE | 820 | THEFT | $500 AND | RESIDENT | FALSE | FALSE | 2532 | 25 | 37 | 25 | 6 | 1140746 | 1909717 | 2012 | 2/10/2018 15:50 | 41.90835 | -87.7584 | (41.908349759, -87.758409329) |
| 8812499 | HV485430 | 9/18/2012 9:00 | 099XX S WALLACE ST | 820 | THEFT | $500 AND | STREET | FALSE | FALSE | 2232 | 22 | 9 | 73 | 6 | 1174056 | 1838887 | 2012 | 2/10/2018 15:50 | 41.71331 | -87.6382 | (41.713307579, -87.638150106) |
| 8812500 | HV485363 | 9/21/2012 11:45 | 027XX N KILPATRICK AVE | 486 | BATTERY | DOMESTIC | APARTME| | TRUE | TRUE | 2521 | 25 | 31 | 19 | 08B | 1144588 | 1917701 | 2012 | 2/10/2018 15:50 | 41.93019 | -87.7441 | (41.930187136, -87.74409405) |
| 8812501 | HV485463 | 9/21/2012 5:00 | 119XX S MICHIGAN AVE | 460 | BATTERY | SIMPLE | STREET | FALSE | FALSE | 532 | 5 | 9 | 53 | 08B | 1178981 | 1825999 | 2012 | 2/10/2018 15:50 | 41.67783 | -87.6205 | (41.677830552, -87.620503762) |
| 8812502 | HV485456 | 9/21/2012 11:30 | 026XX N KOSTNER AVE | 460 | BATTERY | SIMPLE | SIDEWALK | FALSE | FALSE | 2524 | 25 | 31 | 20 | 08B | 1146593 | 1917207 | 2012 | 2/10/2018 15:50 | 41.92879 | -87.7367 | (41.928793506, -87.736738726) |
| 8812503 | HV485362 | 9/7/2012 0:00 | 014XX W 63RD ST | 820 | THEFT | $500 AND | OTHER | FALSE | FALSE | 713 | 7 | 16 | 67 | 6 | 1167639 | 1862965 | 2012 | 2/10/2018 15:50 | 41.77952 | -87.661 | (41.779520772, -87.660962362) |
| 8812504 | HV485233 | 9/21/2012 9:45 | 067XX S LAFAYETTE AVE | 610 | BURGLAR\ | FORCIBLE | RESIDENC | FALSE | FALSE | 722 | 7 | 6 | 69 | 5 | 1176996 | 1860224 | 2012 | 2/10/2018 15:50 | 41.77179 | -87.6267 | (41.771793286, -87.626741205) |
| 8812506 | HV480859 | 9/18/2012 10:25 | 013XX S LAWNDALE AVE | 2014 | NARCOTIC | MANU/DE | SIDEWALK | TRUE | FALSE | 1011 | 10 | 24 | 29 | 18 | 1151920 | 1893559 | 2012 | 2/10/2018 15:50 | 41.8638 | -87.7178 | (41.863797669, -87.717787055) |
| 8812509 | HV485138 | 9/20/2012 15:00 | 063XX S HOYNE AVE | 810 | THEFT | OVER $50( | STREET | FALSE | FALSE | 726 | 7 | 15 | 67 | 6 | 1163449 | 1862682 | 2012 | 2/10/2018 15:50 | 41.77883 | -87.6763 | (41.778833058, -87.676331372) |

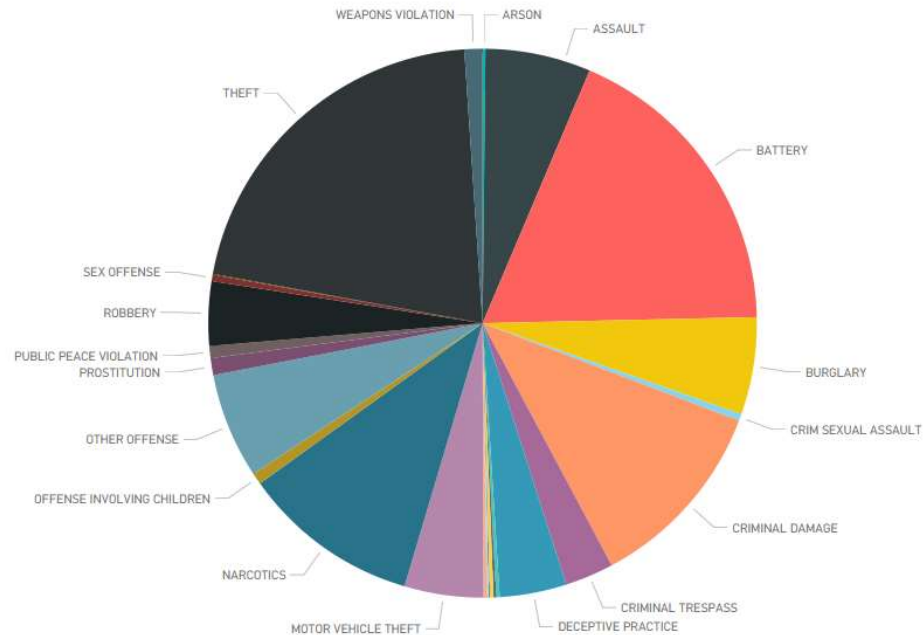I have used MapReduce, Hive, Pig and with various algorithm as stated below:

a. Counter
b. Average
c. Percentage
d. Inverted Index
e. Secondary Sorting
f. Binning
g. Partitioning Pruning
h. Job Chaining
i. Top K filter
j. Bloom Filter
**k.** Normal Regex Filter

# Data Analysis using PowerBi:

Please find all the implementation code in appendix section.
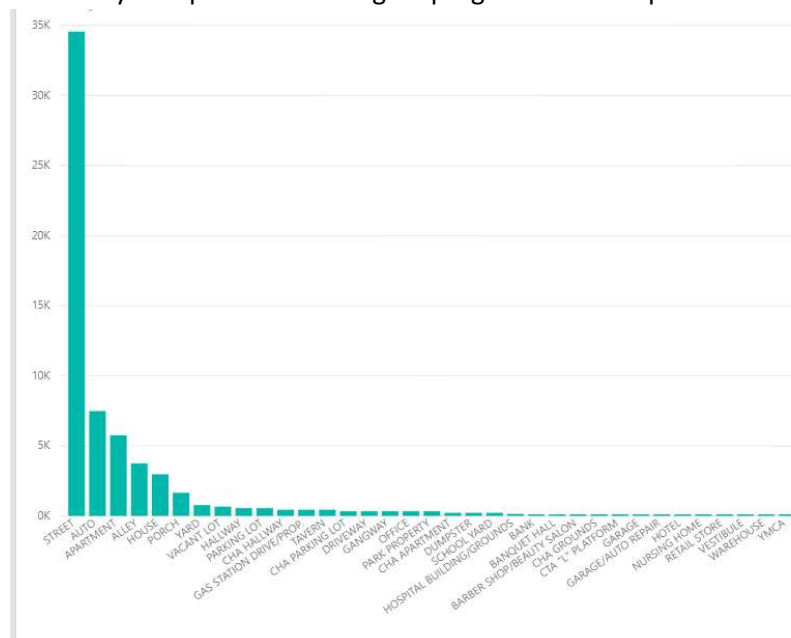
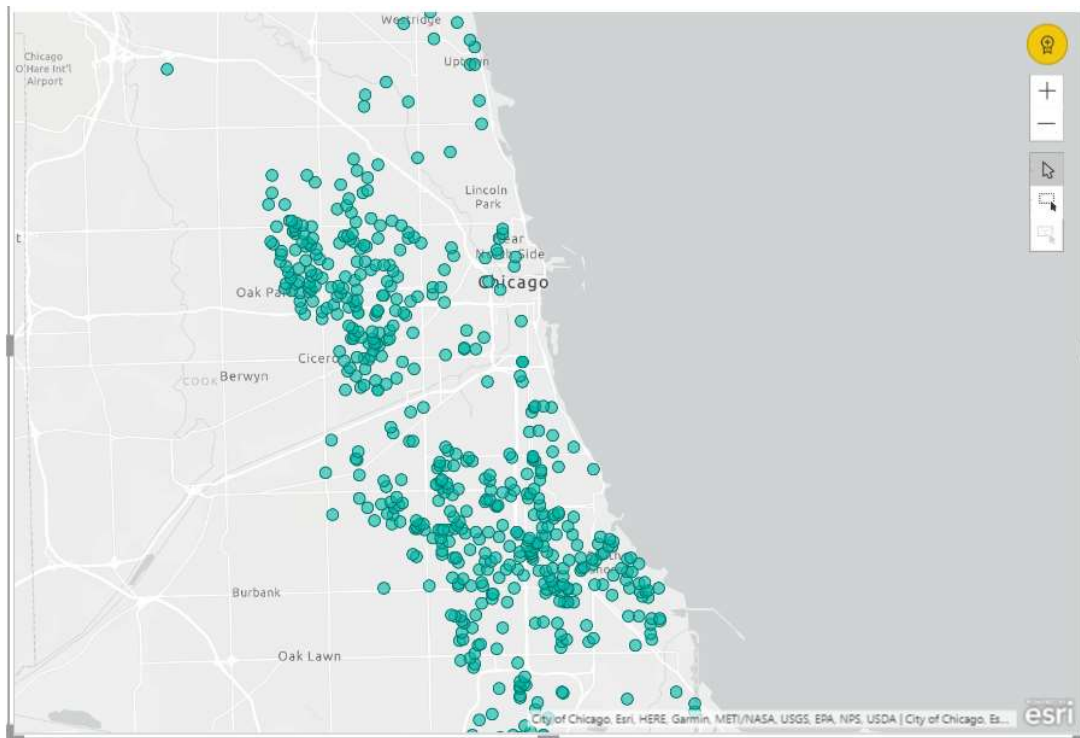1. How many different types of crimes are reported?
   This analysis has been achieved through summarization algorithm in mapreduce.



2. Filters based on crime: Analysis of crime Homicide in different areas of Chicago:
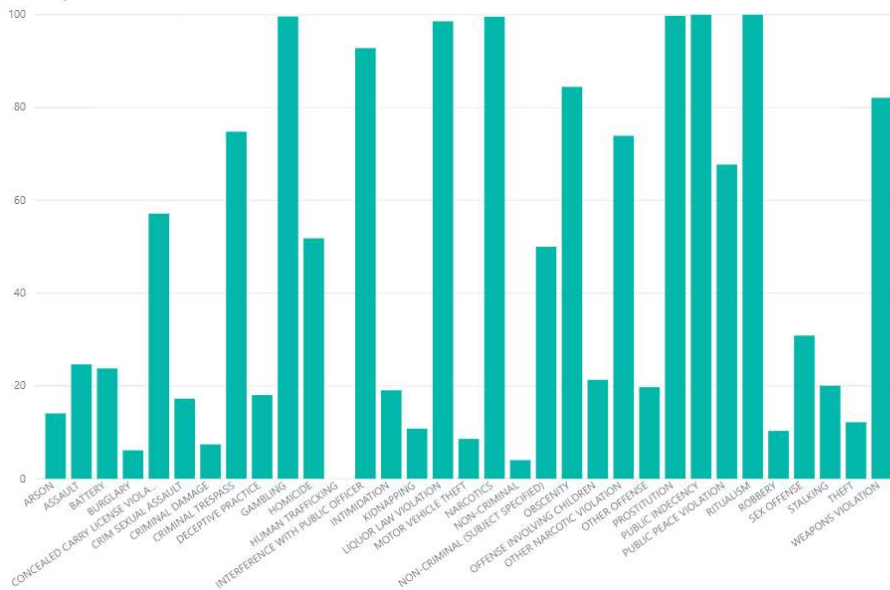   This analysis is performed using mapregex filter in mapreduce

3. Top 25 blocks that has highest number of crime incidents.
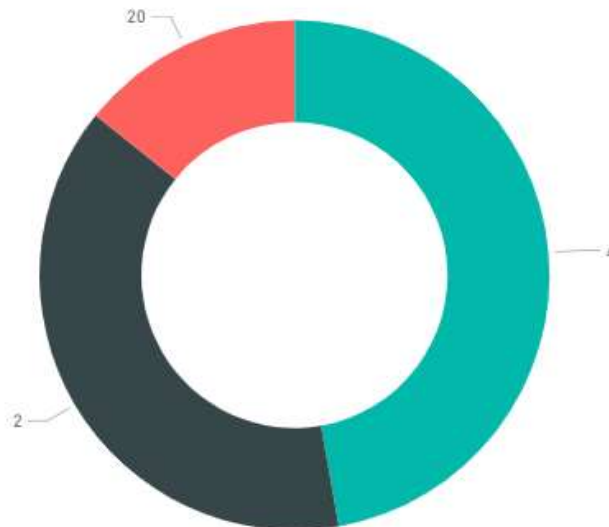   This is achieved through pig script.

| | |
|---|---|
| 001XX N STATE ST | 2243 |
| 076XX S CICERO AVE | 1607 |
| 100XX W OHARE ST | 1560 |
| 008XX N MICHIGAN AVE | 1446 |
| 0000X W TERMINAL ST | 1329 |
| 0000X N STATE ST | 1239 |
| 009XX W BELMONT AVE | 771 |
| 064XX S DR MARTIN LUTHER KING JR DR | 746 |
| 008XX N STATE ST | 744 |
| 051XX W MADISON ST | 712 |
| 063XX S DR MARTIN LUTHER KING JR DR | 706 |
| 011XX W WILSON AVE | 685 |
| 001XX W 87TH ST | 636 |
| 040XX W LAKE ST | 629 |
| 083XX S STEWART AVE | 611 |
| 0000X S STATE ST | 602 |
| 001XX W LAKE ST | 580 |
| 002XX W 87TH ST | 577 |
| 046XX N BROADWAY | 565 |
| 006XX N MICHIGAN AVE | 551 |
| 012XX S WABASH AVE | 547 |
| 007XX N MICHIGAN AVE | 541 |
| 062XX S DR MARTIN LUTHER KING JR DR | 529 |
| 042XX W MADISON ST | 528 |
| 046XX W NORTH AVE | 524 |
| **Total** | 21208 |

4. Percentage of cases resolved under each crime categories (arrest percentage).
   This is achieved through custom writable class in mapreduce calculating percentage of arrest
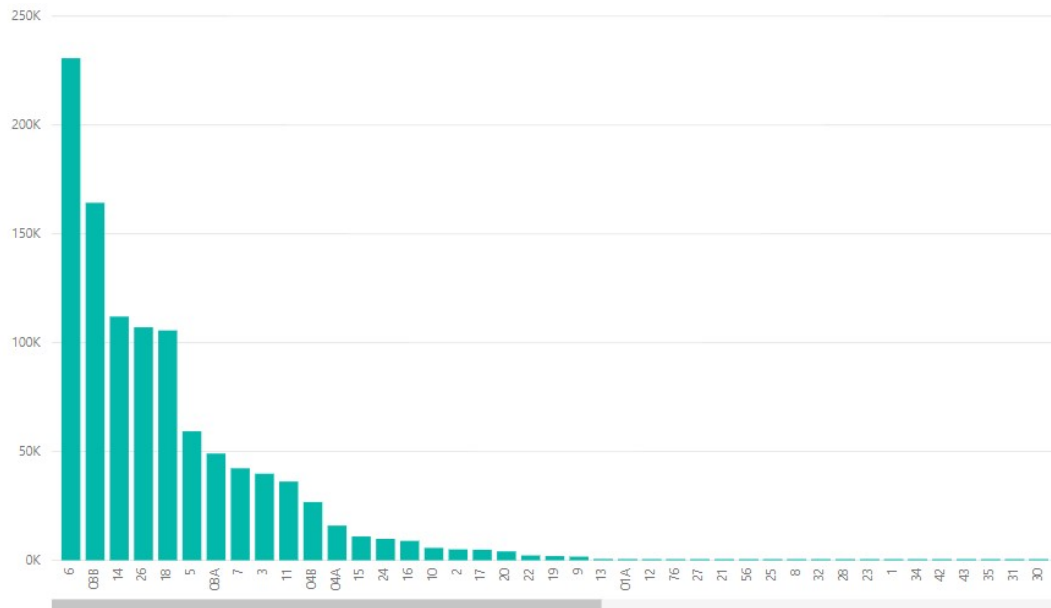


5. Count of crime filtered through district.
   This is achieved using bloom filter where 3 district value has been passed to bloom filter
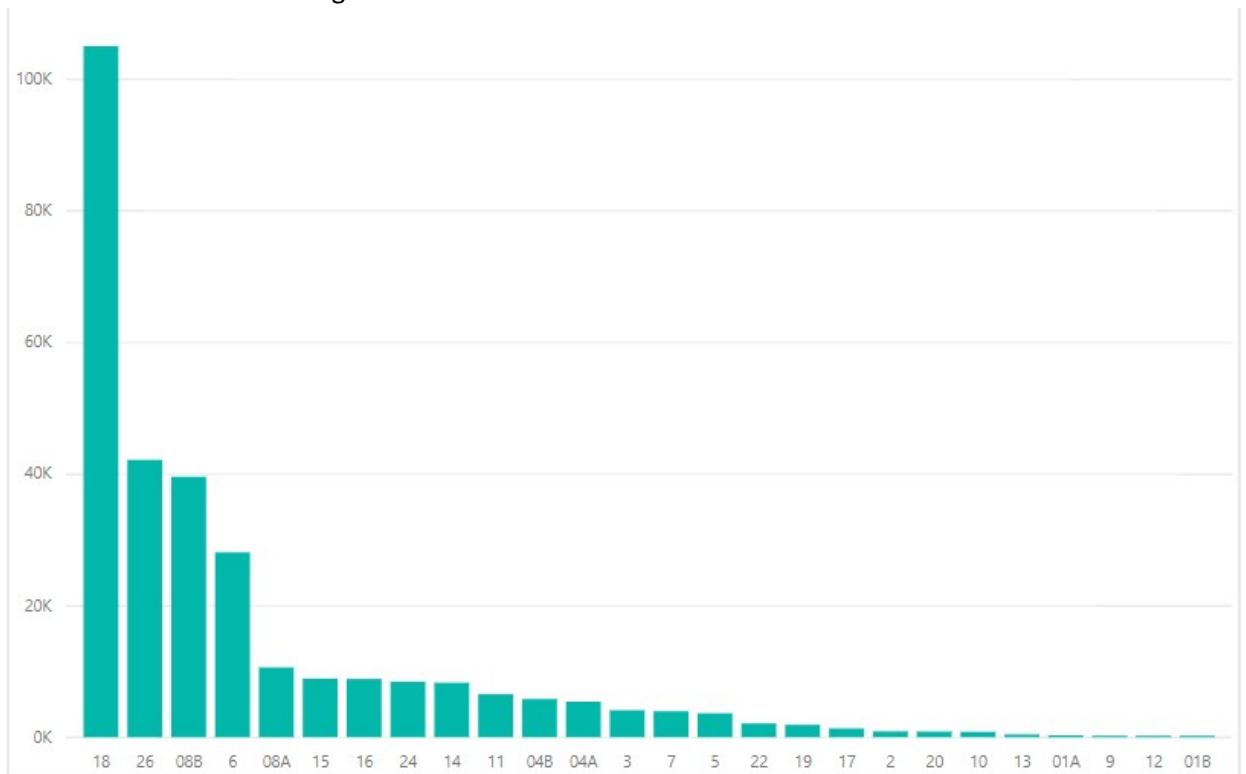
Count of Crime by District

6. Total cases with FBI.
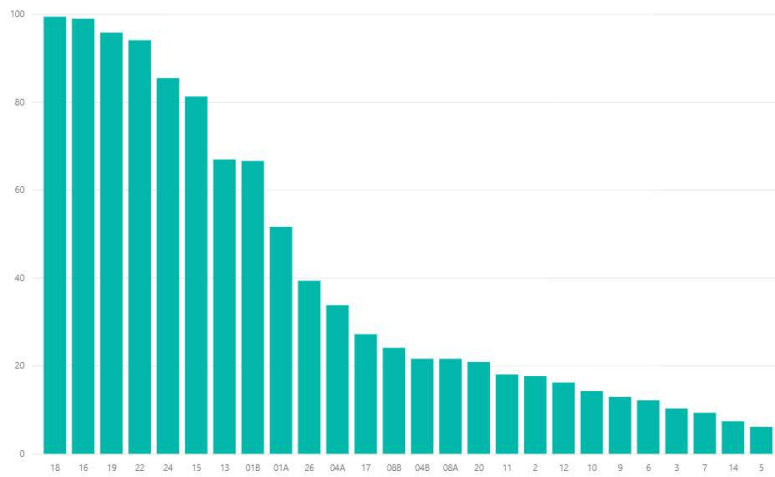   This is achieved in hive where we analyzed total cases with FBI.



7. Total cases resolved by FBI overall
   This is achieved in hive using filter on cases resolved value = true.

8. Cases solved in percentage:
   This analysis uses job chaining algorithm where first job calculates the percentage of arrest and 2nd job sorts the data in descending order.
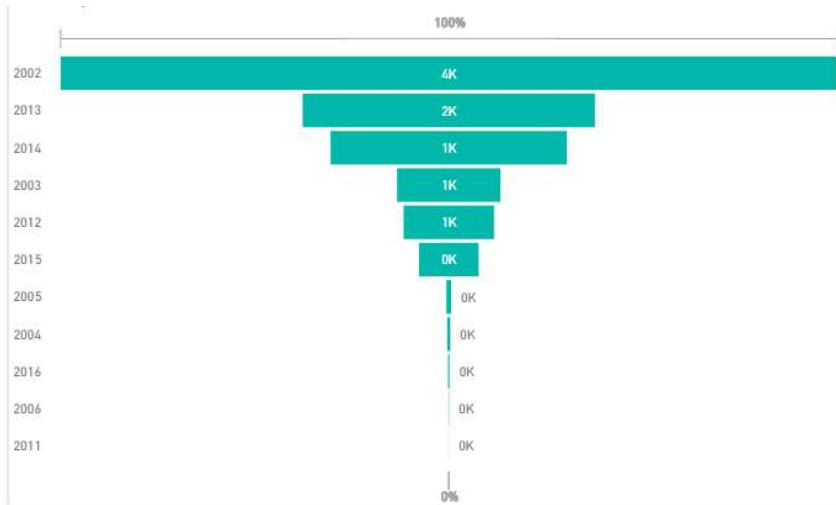


9. Domestic cases which are not resolved
   This is achieved in Hive script using filters on domestic values and arrest values.

| Column1 | Column2 |
|---|---|
| BATTERY | 65224 |
| OTHER OFFENSE | 16782 |
| ASSAULT | 12228 |
| CRIMINAL DAMAGE | 8479 |
| THEFT | 5638 |
| OFFENSE INVOLVING CHILDREN | 2838 |
| CRIMINAL TRESPASS | 613 |
| ROBBERY | 548 |
| DECEPTIVE PRACTICE | 426 |
| CRIM SEXUAL ASSAULT | 417 |
| KIDNAPPING | 313 |
| PUBLIC PEACE VIOLATION | 305 |
| BURGLARY | 281 |
| SEX OFFENSE | 235 |
| MOTOR VEHICLE THEFT | 185 |
| STALKING | 182 |
| INTIMIDATION | 57 |
| ARSON | 51 |
| WEAPONS VIOLATION | 4 |
| OBSCENITY | 3 |
| HOMICIDE | 2 |
| INTERFERENCE WITH PUBLIC OFFICER | 2 |
| NARCOTICS | 1 |
| NON-CRIMINAL (SUBJECT SPECIFIED) | 1 |
| **Total** | **114815** |

10. Crime analysis based on year
    This is achieved in pig through secondary sorting algorithm.
    First data is grouped with crime category and year and then flattening the group parameter to achieve sorting first on crime category and then on year.

Above analysis has been achieved through MapReduce, Hive, Pig and with various algorithm as stated below:

l. Counter
m. Average
n. Percentage
o. Inverted Index
p. Secondary Sorting
q. Binning
r. Partitioning Pruning
s. Job Chaining
t. Top K filter
u. Bloom Filter
v. Normal Regex Filter

I have implemented all the above algorithm using Map-reduce pig and hive.

Apart from these algorithms, mahout has also been used to predict crime that can happen in different districts based on **mahout's recommendation** algorithm framework.

 Below is the snapshot which says district 1 has 99% of chance that its unsafe for Kidnapping and should avoid playing in night in that area.

```
19/04/26 12:30:53 INFO model.GenericDataModel: Processed 23 users
district Id: 1
Chances of crime committed 581. Strength of the preference: 99.000000
Chances of crime committed 1520. Strength of the preference: 37.809219
Chances of crime committed 450. Strength of the preference: 37.247375
Chances of crime committed 4240. Strength of the preference: 37.247375
Chances of crime committed 462. Strength of the preference: 25.721134
district Id: 2
Chances of crime committed 1900. Strength of the preference: 60.008011
Chances of crime committed 275. Strength of the preference: 45.866837
Chances of crime committed 1566. Strength of the preference: 39.001488
Chances of crime committed 1435. Strength of the preference: 31.232704
Chances of crime committed 1255. Strength of the preference: 30.165001
district Id: 3
Chances of crime committed 1631. Strength of the preference: 67.060188
Chances of crime committed 3920. Strength of the preference: 56.016571
```

**Appendix: Implementation Code:**

# Analysis using MapReduce:

1. **Percentage of Arrest:**
   a. Map Class

```java
import java.io.IOException;

public class Map extends Mapper<Object, Text, Text, BooleanWritable> {
    private BooleanWritable arrest1 = new BooleanWritable();
    private Text crime1 = new Text();

    public void map(Object key, Text value, Context context) throws IOException,InterruptedException
    {

        String line = value.toString();
        if(line.toLowerCase().contains("location")) {
            return;
        }
//      System.out.println(line);
        String arrest = line.split(",")[8];
        String crime = line.split(",")[5];
        boolean arrestBool = Boolean.parseBoolean(arrest);
        arrest1.set(arrestBool);
        crime1.set(crime);

        context.write(crime1,arrest1);
    }

}
```

b. Reduce Class

```java
3⊕ import java.io.IOException;□
1
2
3
4 public class Reduce extends Reducer<Text, BooleanWritable, Text, Text> {
5     private Text result = new Text();
6⊖     @Override
7     public void reduce(Text key, Iterable<BooleanWritable> values, Context context)
8             throws IOException, InterruptedException {
9 //        System.out.println("here reducer;");
0         float true_count = 0;
1         float false_count = 0;
2         float count = 0;
3         float percentage_solved = 0;
4         boolean first = true;
5
6         for(BooleanWritable arrest : values) {
7             count++;
8 //            System.out.println(arrest.toString());
9             if(arrest.get() == true) {
0 //                System.out.println(arrest.toString());
1                 true_count++ ;
2 //                System.out.println(true_count);
3             }else {
4 //                sb.append(",");
5                 false_count++;
6 //                System.out.println(false_count);
7             }
8 //            sb.append(crime.toString());
9         }
0 //        System.out.println(true_count);
1
2 //        System.out.println(true_count/count);
3         percentage_solved = (true_count/count)*100;
4 //        System.out.println(percentage_solved);
5         result.set(Float.toString(percentage_solved));
6         context.write(key,result);
7     }
8 }
```

## 2. Binning Class based on different crimes category:

a. Driver class:

```java
📄 App.java ⊠
 1 package Binning;
 2
🔒 3⊕ import java.io.File;□
22
23 public class App {
24⊖     public static void main(String[] args)throws IOException, InterruptedException, ClassNotFoundException {
25         // TODO code application logic here
26         Configuration conf = new Configuration();
27         Job job = Job.getInstance(conf, "Binnin");
28         job.setJarByClass(App.class);
29
30         //Setting Mapper Class and the output key and value
31         job.setMapperClass(BinningMapper.class);
32         job.setMapOutputKeyClass(Text.class);
33         job.setMapOutputValueClass(NullWritable.class);
34
35         //No combiner, partitioner or reducer is used in this pattern!
36         job.setNumReduceTasks(0);
37         FileUtils.deleteDirectory(new File(args[1]));
38         TextInputFormat.addInputPath(job, new Path(args[0]));
39         FileOutputFormat.setOutputPath(job, new Path(args[1]));
40
41         MultipleOutputs.addNamedOutput(job, "bins", TextOutputFormat.class, Text.class, NullWritable.class);
42         MultipleOutputs.setCountersEnabled(job, true);
43
44         System.exit(job.waitForCompletion(true) ? 0 : 1);
45     }
46 }
```

b. Map Class:

```java
1  package Binning;
2
3
4  import java.io.IOException;
19
20 public class BinningMapper extends Mapper<LongWritable, Text, Text, NullWritable> {
21     private MultipleOutputs<Text, NullWritable> mos = null;
22
23         protected void setup(Context context){
24             mos = new MultipleOutputs(context);
25 //          System.out.println("inside");
26 //          System.out.println(mos);
27         }
28
29     protected void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
30
31         if(key.get()==0){
32         return;
33     }
34         String[] token = value.toString().split(",");
35         System.out.println(token);
36             //for(String rate: token){
37
38             String crime_type = token[5].trim();
39             Map map = new HashMap();
40             if (!map.containsKey(crime_type)){
41             map.put(crime_type, crime_type);
42             }
43
44             if(map.containsKey(crime_type)){
45                 mos.write("bins", value, NullWritable.get(), map.get(crime_type).toString());
46             }
47 //
48 }
49         protected void cleanup(Context context) throws IOException, InterruptedException{
50             mos.close();
51         }
52
53
54 }
```

**3. Count of different crime categories:**

a. Map Class

```java
1  package Final_project;
2
3  import java.io.IOException;
13
14 public class Map extends MapReduceBase implements Mapper{
15
16     private final static IntWritable one = new IntWritable(1);
17     private Text word = new Text();
18
19     public void map(Object key, Object value, OutputCollector output, Reporter reporter) throws IOException {
20 //
21         String line = value.toString();
22 //
23         if(line.toLowerCase().contains("location")) {
24             return;
25         }
26         String IP = line.split(",")[5].trim();
27         System.out.println(IP);
28 //
29         word.set(IP);
30         output.collect(word, one);
31
32     }
33 }
34
```

b. Reduce Class:

```java
Reduce.java
1  package Final_project;
2
3  import java.io.IOException;
11
12 public class Reduce extends MapReduceBase implements Reducer {
13
14     public void reduce(Object key, Iterator values, OutputCollector output, Reporter reporter) throws IOException {
15         // TODO Auto-generated method stub
16         int sum = 0;
17         while (values.hasNext()) {
18                 sum += ((IntWritable) values.next()).get();
19         }
20         output.collect(key, new IntWritable(sum));
21     }
22 }
23 //
```

## 4. Inverted Index for street vs location search

### a. Map Class:

```java
import java.io.IOException;

public class Map extends Mapper<Object, Text, Text, Text> {
    private Text blockadrr = new Text();
    private Text crimetype = new Text();

    public void map(Object key, Text value, Context context) throws IOException,InterruptedException
    {

        String line = value.toString();
//      System.out.println(line);
        String IP = line.split(",")[3];
        String urls = line.split(",")[7];

//      System.out.println(urls);
        blockadrr.set(IP);
        crimetype.set(urls);

        context.write(blockadrr, crimetype);
    }

}
```

### b. Reducer class:

```java
package InvertedIndex;

import java.io.IOException;


public class Reduce extends Reducer<Text, Text, Text, Text> {
    private Text result = new Text();
    @Override
    public void reduce(Text key, Iterable<Text> values, Context context)
            throws IOException, InterruptedException {
        try {
        HashMap m=new HashMap();
        int count=0;
        StringBuilder sb = new StringBuilder();
        boolean first = true;

        for(Text crime_Type : values) {
            String str=crime_Type.toString();
            if(m!=null &&m.get(str)!=null){
                count+= 1;
                m.put(str, ++count);
            }else{
                /*Else part will execute if file name is already added then just
                m.put(str, 1);
            }
//          if(first) {
//              first=false;
//          }else {
//              sb.append(",");
//          }
//          sb.append(crime_Type.toString());
        }

//      result.set(sb.toString());
//      context.write(key,result);
        context.write(key, new Text(m.toString()));
        }
        catch(Exception e) {
            System.out.println(e);
        }
    }
}
```

## 5.Job Chaining:

```java
 2⊕ * To change this license header, choose License Headers in Project Properties.
 6  package JobChaining;
 7
 8⊕ import org.apache.commons.io.FileUtils;
23
24  public class App {
25
26⊖     /**
27       * @param args the command line arguments
28       */
29⊖     public static void main(String[] args) throws IOException, InterruptedException, ClassNotFoundException{
30          Configuration conf1 = new Configuration();
31          Job job1 = Job.getInstance(conf1,"Amazon Average");
32          job1.setJarByClass(App.class);
33          job1.setMapperClass(Map1.class);
34          job1.setMapOutputKeyClass(Text.class);
35          job1.setMapOutputValueClass(BooleanWritable.class);
36
37          job1.setReducerClass(Reduce1.class);
38          job1.setOutputKeyClass(Text.class);
39          job1.setOutputValueClass(FloatWritable.class);
40          FileUtils.deleteDirectory(new File(args[1]));
41          FileUtils.deleteDirectory(new File(args[2]));
42          FileInputFormat.addInputPath(job1, new Path(args[0]));
43          //FileOutputFormat.setOutputPath(job1, new Path(args[1]));
44          FileOutputFormat.setOutputPath(job1, new Path(args[1]));
45          boolean complete = job1.waitForCompletion(true);
46
47          Configuration conf2 = new Configuration();
48          Job job2 = Job.getInstance(conf2,"Chaining Sorting");
49
50          if(complete){
51          job2.setJarByClass(App.class);
52          job2.setMapperClass(Map2.class);
53          job2.setMapOutputKeyClass(FloatWritable.class);
54          job2.setMapOutputValueClass(Text.class);
55
56          job2.setReducerClass(Reduce2.class);
57          job2.setOutputKeyClass(Text.class);
58          job2.setOutputValueClass(FloatWritable.class);
59
60          FileInputFormat.addInputPath(job2, new Path(args[1]));
61          FileOutputFormat.setOutputPath(job2, new Path(args[2]));
62
63          System.exit(job2.waitForCompletion(true) ? 0:1);
64
65          }
66      }
67
68⊖     public static class Map1 extends Mapper<LongWritable, Text, Text, BooleanWritable>{
69
70          private Text text = new Text();
71          private BooleanWritable arrested = new BooleanWritable();
72
73⊖         protected void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
74
75              if(key.get()==0){
76              return;
77              }
78
79              else{
80
81              String[] line = value.toString().split(",");
82              if (line[8]=="arrest") {
83                  return;
84              }
85              String fbiCode = line[14].trim();
86              boolean arrest = Boolean.valueOf(line[8].trim());
87              System.out.println(arrest);
88
89              text.set(fbiCode);
90              arrested.set(arrest);
91
92              context.write(text, arrested);
93
94              }
95          }
96      }
97
98⊖      public static class Reduce1 extends Reducer<Text, BooleanWritable, Text, FloatWritable>{
99
100         private FloatWritable result = new FloatWritable();
101
102⊖        @Override
103         protected void reduce(Text key, Iterable<BooleanWritable> values, Context context) throws IOException, InterruptedException{
104
105             float sum = 0;
106             int count = 0;
```

```
            for(BooleanWritable val:values){
                System.out.println(val.get());
                if (val.get() == true) {
                    sum += 1;
//                  System.out.println("float");
                }
                count = count+1;
            }

            float average = (sum/count)*100;

            result.set(average);
            context.write(key,result);
        }

    }

    public static class Map2 extends Mapper<LongWritable, Text, FloatWritable, Text>{

        public void map(LongWritable key, Text value, Context context){

            String[] row =value.toString().split("\\t");
            Text fbi = new Text(row[0]);
            float counting = Float.valueOf(row[1].trim());

            try{
                FloatWritable count = new FloatWritable(counting);
                context.write(count, fbi);
            }

            catch(Exception e){

            }
        }
    }

    public static class Reduce2 extends Reducer<FloatWritable, Text, Text, FloatWritable>{

        public void reduce(FloatWritable key, Iterable<Text> value, Context context) throws IOException, InterruptedException{

            for(Text val : value){
```

```
150         context.write(val,key);
151     }
152   }
153   }
154 }
155
```

## 6.Bloom Filter:

### a. Driver Class

```
 1 package LocationFilter;
 2
 3⊕ import org.apache.hadoop.fs.FileSystem;
10
11 public class Driver {
12
13⊖    public static void main(String[] args) throws Exception {
14
15         Path inputPath = new Path(args[0]);
16         Path outputDir = new Path(args[1]);
17         Job job = Job.getInstance();
18         job.setJarByClass(Driver.class);
19         job.setMapperClass(BloomFilterMapper.class);
20         job.setMapOutputKeyClass(Text.class);
21         job.setMapOutputValueClass(NullWritable.class);
22         job.setNumReduceTasks(0);
23         FileInputFormat.addInputPath(job, inputPath);
24         FileOutputFormat.setOutputPath(job, outputDir);
25         FileSystem hdfs = FileSystem.get(job.getConfiguration());
26         if (hdfs.exists(outputDir))
27             hdfs.delete(outputDir, true);
28
29         job.waitForCompletion(true);
30
31     }
32 }
33
```

b. Location pojo filter:

```
1  package LocationFilter;
2
3  public class Location {
4
5      final String district;
6
7
8      Location(String district){
9
10         this.district = district;
11
12
13     }
14 }
```

c. Bloom Filter:

```
1  package LocationFilter;
2
3  import java.io.IOException;
15
16 public class BloomFilterMapper extends Mapper<LongWritable, Text, Text, NullWritable> {
17
18     private BloomFilter<Location> friends;
19     Funnel<Location> p = new Funnel<Location>() {
20         public void funnel(Location from, Sink into) {
21             into.putString(from.district, Charsets.UTF_8);
22         }
23     };
24     @Override
25     protected void setup(Mapper<LongWritable, Text, Text, NullWritable>.Context context)
26             throws IOException, InterruptedException {
27         this.friends = BloomFilter.create(p, 500, 0.1);
28         Location p1 = new Location("20");
29         Location p2 = new Location("4");
30         Location p3 = new Location("2");
31         ArrayList<Location> friendList = new ArrayList<Location>();
32         friendList.add(p1);
33         friendList.add(p2);
34         friendList.add(p3);
35         for (Location p : friendList) {
36             friends.put(p);
37         }
38
39     }
40
41     @Override
42     protected void map(LongWritable key, Text value, Mapper<LongWritable, Text, Text, NullWritable>.Context context)
43             throws IOException, InterruptedException {
44
45         String values[] = value.toString().split(",");
46         Location p = new Location(values[11]);
47
48         if (friends.mightContain(p)) {
49             context.write(value, NullWritable.get());
50         }
51     }
52 }
53
```

### 7. Regex Filter to filter crime category:

Driver Class:

```
package NormalFilter;

 * Hello world!
import java.io.IOException;

public class App
{
    public static void main( String[] args ) throws IOException, ClassNotFoundException, InterruptedException
    {
        long start_time = System.currentTimeMillis();
        Configuration conf = new Configuration();
        conf.set("mapregex", "HOMICIDE");
        Job job = Job.getInstance(conf, "Stocks");
        job.setJarByClass(App.class);
        job.setMapperClass(Mapp.class);
        job.setMapOutputKeyClass(NullWritable.class);
        job.setMapOutputValueClass(Text.class);
//          job.setReducerClass(Reduce.class);
        job.setNumReduceTasks(1);
//          job.setOutputKeyClass(Text.class);
//          job.setOutputValueClass(BooleanWritable.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        Path outDir = new Path(args[1]);
        FileSystem fs = FileSystem.get(job.getConfiguration());
        if(fs.exists(outDir)) {
            fs.delete(outDir, true);
        }

        System.exit(job.waitForCompletion(true) ? 0 : 1);
        System.out.println("Time Taken is :" + (System.currentTimeMillis() - start_time)/ 1000);
    }
}
```

Map Class:

```
1  package NormalFilter;
2
3
4⊕ import java.io.DataInputStream;
12
13  public class Mapp extends Mapper<Object, Text, NullWritable, Text> {
14
15      private String mapRegex = null;
16
17⊝     public void setup(Context context) throws IOException,
18          InterruptedException {
19
20          mapRegex = context.getConfiguration().get("mapregex");
21  //      System.out.println(mapRegex);
22      }
23
24⊝     public void map(Object key, Text value, Context context)
25          throws IOException, InterruptedException {
26  //      System.out.println(value.toString().split(",")[7]);
27          String location = value.toString().split(",")[5];
28          if (location.matches(mapRegex)) {
29
30              context.write(NullWritable.get(), value);
31          }
32      }
33  }
```

## 8.Partition Pruning to partition data based on years:

Driver Class:

```
App.java ⊠
1   package partition_pruning;
2
4⊕  * Hello world!
7⊕ import java.io.File;
15  |
16  public class App {
17⊝ public static void main(String[] args) throws Exception {
18          FileUtils.deleteDirectory(new File(args[1]));
19          if (args.length != 2) {
20          System.err.println("Please specify the input and output path");
21          System.exit(-1);
22          }
23          Configuration conf = new Configuration();
24          Job job = Job.getInstance(conf);
25          job.setPartitionerClass(CustomPartitioner.class);
26          job.setJarByClass(App.class);
27          job.setJobName("Partioning_Pattern");
28          FileInputFormat.addInputPath(job, new Path(args[0]));
29          FileOutputFormat.setOutputPath(job, new Path(args[1]));
30          job.setMapperClass(Map.class);
31          job.setReducerClass(Reduce.class);
32          job.setOutputKeyClass(Text.class);
33          job.setOutputValueClass(Text.class);
34          job.setNumReduceTasks(18);
35          System.exit(job.waitForCompletion(true) ? 0 : 1);
36  }
37  }
```

Map Class:

```
package partition_pruning;

⊕ import java.io.IOException;

public class Map extends Mapper<Object, Text, Text, Text> {
    public Text mapperKey=new Text();
⊝   @Override
    public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
    String data = value.toString();
    String[] field = data.split(",");
    if (field[17].trim().length() > 0 && field[17].trim().length() < 5) {
    String departmentName = field[17];
    mapperKey.set(departmentName);
    context.write(mapperKey, value);
    }
}
}
```

Custom Partitioner Class:

```
1  package partition_pruning;
2
3⊕ import org.apache.hadoop.conf.Configurable;▯
7
8  public class CustomPartitioner extends Partitioner<Text, Text>implements Configurable {
9      private Configuration conf = null;
0⊖     public void setConf(Configuration conf) {
1      this.conf = conf;
2      }
3  //  @Override
4⊖     public Configuration getConf() {
5      return conf;
6      }
7  //  @Override
8⊖     public int getPartition(Text key, Text value, int numPartitions) {
9      return Math.abs((key.toString().hashCode()) % numPartitions);
0      }
1  }
```

Reducer Class:

```
1  package partition_pruning;
2
3⊕ import java.io.IOException;▯
7
3  public class Reduce extends Reducer<Text, Text, Text, NullWritable> {
3⊖ @Override
9      public void reduce(Text text, Iterable<Text> values, Context context) throws IOException, InterruptedException {
1      for (Text value : values) {
2      context.write(value, NullWritable.get());
3      }
4      }
5  }
```

# Analysis through Pig Script:

1. Count of cases with FBI:

```
crime_data = LOAD '/home/bhavik/Desktop/bigdata/Crimes_cleaned.csv' USING PigStorage(',')
as (ID:chararray,Case_Number:chararray,Date:chararray,Block:chararray,
IUCR:chararray,Primary_Type:chararray,Description:chararray,Location_Description:chararray,Arrest:Boolean,
Domestic:Boolean,Beat:chararray,District:chararray,Ward:chararray,
Community_Area:chararray,FBI_Code:chararray,Year:chararray);
rh = filter crime_data by $2 != 'Date';
gr = GROUP rh BY (FBI_Code);
gen = FOREACH gr GENERATE group AS (Fbi), COUNT(rh.FBI_Code);
orders = ORDER gen BY $1 DESC;
STORE orders into 'fbi' using PigStorage(',');
```

2. Secondary sorting separating data based on year and categories of crime:

```
crime_data = LOAD '/home/bhavik/Desktop/bigdata/Crimes_cleaned.csv' USING PigStorage(',')
as (ID:chararray,Case_Number:chararray,Date:chararray,Block:chararray,
IUCR:chararray,Primary_Type:chararray,Description:chararray,Location_Description:chararray,
Arrest:Boolean,Domestic:Boolean,Beat:chararray,District:chararray,Ward:chararray,
Community_Area:chararray,FBI_Code:chararray,Year:chararray);
remove_header = filter crime_data by $2 != 'Date';
gr = GROUP remove_header BY (Primary_Type,Year);
gen = FOREACH gr GENERATE FLATTEN(group) AS (crime, year), COUNT(remove_header.Primary_Type) as crime_count;
fil = filter gen by $0 == 'THEFT';
orders = ORDER gen BY $0 ASC, $1 DESC;
STORE orders into 'secondary_sort' using PigStorage(',');
```

3. Top 25 blocks where crimes are reported:

```
crime_data = LOAD '/home/bhavik/Desktop/bigdata/Crimes.csv' USING PigStorage(',') as (ID:chararray,Case_Number:chararray,
Date:chararray,Block:chararray,IUCR:chararray,Primary_Type:chararray,
Description:chararray,Location_Description:chararray,Arrest:Boolean,Domestic:Boolean,Beat:chararray,
District:chararray,Ward:chararray,Community_Area:chararray,FBI_Code:chararray,X_Coordinate::chararray,
Y_Coordinate::chararray,Year:chararray,Updated_On:chararray,Latitude:chararray,Longitude:chararray,Location:chararray);
remove_header = filter crime_data by $2 != 'Date';
gr = GROUP remove_header BY Block;
gen = FOREACH gr GENERATE group as Block, COUNT(remove_header.Block);
orders = ORDER gen BY $1 DESC;
top_location = LIMIT orders 25;
STORE top_location into 'top25_final_block' using PigStorage(',');
```

# Analysis through Hive

1. Percentage of crime solved by FBI:

```
    > INSERT OVERWRITE DIRECTORY '/casesResolved' select FBI_Code, Count(Arrest) from crime
where Arrest == true
group by FBI_Code;
```

2. Domestic cases which are not resolved:

```
hive> INSERT OVERWRITE DIRECTORY '/domestic' select Primary_Type, count(Primary_Type) from crime
    > where Arrest == false
    > and Domestic == true
 group by Primary_Type;
```

# Mahout Recommendation Algorithm

Algorithm to predict crime that can happen in different districts:

```java
package safetyzone;

import java.io.File;

public class App
{
    public static void main(String[] args) throws IOException, TasteException{

        File userPreferencesFile = new File("/home/bhavik/Desktop/bigdata/crimes_recommendation.csv");
        DataModel dataModel = new FileDataModel(userPreferencesFile);

        UserSimilarity userSimilarity = new PearsonCorrelationSimilarity(dataModel);
        UserNeighborhood userNeighborhood = new ThresholdUserNeighborhood(0.1, userSimilarity, dataModel);

        // Create a generic user based recommender with the dataModel, the userNeighborhood and the userSimilarity
        Recommender genericRecommender = new GenericUserBasedRecommender(dataModel, userNeighborhood, userSimilarity);

        // Recommend 5 items for each user
        for (LongPrimitiveIterator iterator = dataModel.getUserIDs(); iterator.hasNext();)
        {
            long district = iterator.nextLong();

            // Generate a list of 5 recommendations for the user
            List<RecommendedItem> itemRecommendations = genericRecommender.recommend(district, 5);

            System.out.format("district Id: %d%n", district);

            if (itemRecommendations.isEmpty())
            {
                System.out.println("safest district.");
            }
            else
            {

                for (RecommendedItem recommendedItem : itemRecommendations)
                {
                    System.out.format("Chances of crime committed %d. Strength of the preference: %f%n", recommendedItem.getItemID(), recommendedItem.getValue());
                }
            }
        }
    }
}
```