# Big Data Crime Analysis Lab Tutorial

**Authors: Jesus Perez, Abhinav Singh, Tien Cheng**
**Instructor: Jongwook Woo**

# Lab Tutorial: Analyzing Crime Data with HiveQL

**Objectives:** This tutorial guides you through downloading, uploading, processing, and analyzing crime datasets from Los Angeles (2010–Present) and Chicago (2001–Present) using Hadoop HDFS and HiveQL. You'll create tables, query the data, and export results as CSV files for visualization. Follow the steps carefully to ensure a smooth workflow.

**Platform Spec:** Oracle Linux Server v. 7.9, 5 Nodes (2 Master, 3 Worker), 31GBx5 RAM, 2.5GHz CPU for smooth data handling.

## Prerequisites

- Access to Oracle Linux server, Hadoop cluster with HDFS tool and Hive installed to run beeline.

- SSH client (e.g., Git Bash, PuTTY) for file transfer and cluster access.

- Basic familiarity with SQL and command-line interfaces.

- SCP tool for file transfers (available in Git Bash or similar terminals).

## Step 1: Download the Datasets

Download the following datasets and save them to your local machine (e.g., `C:\Users\YourName\Downloads`):

1. **LA Crime Data (2020–Present)**:
   https://data.lacity.org/Public-Safety/Crime-Data-from-2020-to-Present/2nrs-mtv8/about_data
2. **LA Crime Data (2010–2019)**:

3. **Chicago Crime Data (2001–Present)**:
   [https://data.cityofchicago.org/Public-Safety/Crimes-2025/t7ek-mgzi/about_data](https://data.cityofchicago.org/Public-Safety/Crimes-2025/t7ek-mgzi/about_data)

# Step 2: Upload Datasets to HDFS

## 2.1 Transfer Files to Cluster

Use SCP to transfer files to the `/tmp` directory of your Hadoop cluster:

```
scp "C:\Users\YourName\Downloads\Crime_Data_from_2010_to_2019.csv"
your_username@144.24.13.0:/tmp/
scp "C:\Users\YourName\Downloads\Crime_Data_from_2020_to_Present.csv"
your_username@144.24.13.0:/tmp/
scp "C:\Users\YourName\Downloads\Crimes_-_2001_to_Present.csv"
your_username@144.24.13.0:/tmp/
```

## 2.2 Verify Upload

```
ssh your_username@144.24.13.0
cd /tmp
ls
```

Ensure you see the uploaded files listed.

# Step 3: Store Files in HDFS

## 3.1 Los Angeles Data

```
hdfs dfs -mkdir crime_data_LA
hdfs dfs -put /tmp/Crime_Data_from_2010_to_2019.csv crime_data_LA
hdfs dfs -put /tmp/Crime_Data_from_2020_to_Present.csv crime_data_LA
```

## 3.2 Chicago Data

```
hdfs dfs -mkdir crime_data_CH
hdfs dfs -put /tmp/Crimes_-_2001_to_Present.csv crime_data_CH
```

# Step 4: Create Hive Tables

**Start Beeline and Connect**

```
> beeline
```

**4.1 Create Los Angeles Table**

```
USE your_username;
```

```
DROP TABLE IF EXISTS crime_data;

CREATE TABLE crime_data (
    DR_NO STRING,
    Date_Rptd STRING,
    DATE_OCC STRING,
    TIME_OCC INT,
    AREA INT,
    AREA_NAME STRING,
    Rpt_Dist_No INT,
    Part_1_2 INT,
    Crm_Cd INT,
    Crm_Cd_Desc STRING,
    Mocodes STRING,
    Vict_Age INT,
    Vict_Sex STRING,
    Vict_Descent STRING,
    Premis_Cd INT,
    Premis_Desc STRING,
    Weapon_Used_Cd STRING,
    Weapon_Desc STRING,
    Status STRING,
    Status_Desc STRING,
    Crm_Cd_1 INT,
    Crm_Cd_2 INT,
```

```
    Crm_Cd_3 INT,
    Crm_Cd_4 INT,
    LOCATION STRING,
    Cross_Street STRING,
    LAT DOUBLE,
    LON DOUBLE
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE
LOCATION '/user/your_username/crime_data_LA'
TBLPROPERTIES ('skip.header.line.count'='1');
```

```
SELECT * FROM crime_data LIMIT 10;
```

## 4.2 Create Chicago Table

```
DROP TABLE IF EXISTS crime_data_CH;

CREATE EXTERNAL TABLE crime_data_CH (
    ID BIGINT,
    Case_Number STRING,
    Dates STRING,
    Block STRING,
    IUCR STRING,
    Primary_Type STRING,
    Description STRING,
    Location_Description STRING,
    Arrest BOOLEAN,
    Domestic BOOLEAN,
    Beat INT,
    District INT,
    Ward INT,
    Community_Area INT,
    FBI_Code STRING,
    X_Coordinate DOUBLE,
    Y_Coordinate DOUBLE,
    Year INT,
    Updated_On STRING,
    Latitude DOUBLE,
    Longitude DOUBLE,
```

```
    Location STRING
)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE
LOCATION '/user/your_username/crime_data_CH'
TBLPROPERTIES ('skip.header.line.count'='1');
```

```
SELECT * FROM crime_data_CH LIMIT 10;
```

# Step 5: Alternative Export via Hive Table

```
NOTE: CSV files don't have the headers when exported so give these headers
to the field in this sequence
```

## 6.1 Create Export Table

```
DROP TABLE IF EXISTS table_csv_export_data;

CREATE TABLE table_csv_export_data
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
LINES TERMINATED BY '\n'
STORED AS TEXTFILE AS
<your entire Hive QL Query here>
```

```
INSERT INTO table_csv_export_data
<your entire Hive QL Query here>
```

## 6.2 Find Table Location

```
DESCRIBE FORMATTED table_csv_export_data;
```

Note the `Location` field in the output.

## 6.3 Extract CSV from HDFS (in a new terminal)

```
ssh your_username@144.24.13.0
hadoop fs -cat /path/to/table_csv_export_data/* > ~/name_your_csv_file.csv
```

## 6.4 Download to Local Machine (in a new terminal)

```
scp your_username@144.24.13.0:/home/your_username/name_your_csv_file.csv
~/name_your_csv_file.csv
```

# Step 6: Analyze the Data (These visualizations are done in Excel 3D Map and Tableau, but you can also do them in Power BI)

## 5.1 Crime Location and Time

**Chicago**

```sql
SELECT
    Dates AS crime_datetime,
    Location,
    Latitude,
    Longitude,
    Primary_Type
FROM crime_data_ch
WHERE Dates IS NOT NULL
  AND Location IS NOT NULL
ORDER BY Dates
LIMIT 5;
```

## Visualize on Excel 3D Map:

**Open Excel** and load your spreadsheet.

Click on **search** > type **3D Map** > **Open 3D Maps**.

A new window will open with a globe or map. Click **New Tour** if prompted.

Excel will auto-detect the location columns (e.g., latitude/longitude or addresses). If not:

- Use the right-side panel to manually assign fields (set `Latitude` and `Longitude` or `City/State` if that's what you have).

**Choose "Primary Type" (count - Not Blank) under the Size selection.**

**Choose the "Primary Type" column to the Category section** of the Layer Pane.

- This will allow different crimes to show in **distinct colors**, just like in our screenshot.

**Drag the Date/Time column to the Time box** (bottom of the Layer Pane) to activate the time slider.

Note: We just visualized for particularly the 4 crimes shown in the screenshot, for both LA and Chicago cities, but one can analyze for whatever number of crimes one wants it to be.



**Los Angeles**

```sql
SELECT
    DATE_OCC AS crime_date,
    TIME_OCC AS crime_time,
    LOCATION,
    LAT AS latitude,
    LON AS longitude,
    Crm_Cd_Desc AS crime_type
```

```
FROM crime_data
WHERE DATE_OCC IS NOT NULL
  AND LOCATION IS NOT NULL
ORDER BY DATE_OCC, TIME_OCC
LIMIT 5;
```



## 5.2 Top Types of Crimes

**Chicago**

```
SELECT
    Primary_Type,
    COUNT(*) AS crime_count
FROM crime_data_ch
GROUP BY Primary_Type
ORDER BY crime_count DESC
LIMIT 10;
```

**Los Angeles**

```sql
SELECT
    Crm_Cd_Desc AS crime_type,
    COUNT(*) AS crime_count
FROM crime_data
GROUP BY Crm_Cd_Desc
ORDER BY crime_count DESC
LIMIT 10;
```

## Visualize on Tableau for Top Types of Crimes:

# Step 1: Load Your Data

1. Open **Tableau Desktop** or **Tableau Public**.
2. Click **"Connect" > "To a File" > "Text File"**.
3. Select your dataset: `crime_data_unified.csv` (or your own cleaned file).
4. Click **"Sheet 1"** at the bottom to start building your visualization.

# Step 2: Drag and Drop Fields to Create the Tree Map

## What You Want:

- **Size** of blocks = Number of crimes
- **Color** = City
- **Text** = Crime Type
- **Filters** = City + Year

## Do the Following:

1. **Drag** `crime_type` **to Rows** or the **Label area** in the Marks card.
2. **Drag** `Number of Records` (or use `CNT(*)`) to **Size** on the Marks card.
3. **Drag** `city` **to Color** on the Marks card.
4. **Drag** `crime_type` **to Label** so the text appears inside the blocks.
5. Change the **Marks type** from "Automatic" to **Tree Map**.
6. **Drag** `occurrence_date` **to Filters** pane.
   - Choose "Years"
   - Select the years you want (or all)
7. **Drag** `city` **to Filters** to select between Chicago / LA.

# Step 3: Final Formatting

1. Click on the **"Label"** on the Marks card:
   - Check **"Show mark labels"**.
   - Choose font size 10–12 for clarity.
2. Click on **"Color"** to adjust contrast (optional).
3. Add a **Title** like:
   "Highest Crime Type in Two Cities"



## 5.3 Police Surveillance Analysis

**Chicago**

```sql
SELECT
    District,
    COUNT(*) AS total_crimes,
    SUM(CAST(Arrest AS INT)) AS arrests,
    ROUND(SUM(CAST(Arrest AS INT)) * 100.0 / COUNT(*), 2) AS
arrest_percentage
FROM crime_data_CH
WHERE District IS NOT NULL
GROUP BY District
```

```
ORDER BY arrest_percentage DESC
LIMIT 5;
```

**Los Angeles**

```
SELECT
    AREA_NAME,
    COUNT(*) AS total_crimes,
    SUM(CASE WHEN Status IN ('AA', 'JA') THEN 1 ELSE 0 END) AS arrests,
    ROUND(SUM(CASE WHEN Status IN ('AA', 'JA') THEN 1 ELSE 0 END) * 100.0 /
COUNT(*), 2) AS arrest_percentage
FROM crime_data
WHERE AREA_NAME IS NOT NULL
GROUP BY AREA_NAME
ORDER BY arrest_percentage DESC
LIMIT 5;
```

## Visualize on Tableau for Police Surveillance Analysis:

# Step 1: Load the Data

1. Open **Tableau Desktop** or **Tableau Public**.
2. Click **Connect → Text File** and load your `crime_data_unified.csv`.
3. Go to **Sheet 1**.

# Step 2: Build the Stacked Bar Chart

**Drag and drop the following:**

1. **Drag `location_description`** → to **Columns**
2. **Drag `Number of Records`** (or use `CNT(*)`) → to **Rows**
3. **Drag `city`** → to **Color** on the Marks card
4. **Drag `city`** again → to **Label** (so it shows inside the bar)
5. **Optional:** Drag `Number of Records` again → to **Label** to show the count

# Step 3: Apply Filters

1. **Drag `city`** → to Filters → Select **Chicago and LA**
2. **Drag `occurrence_date`** → to Filters → Choose "Years" → Select years 2010–2025 or as needed
3. Optional: **Drag `location_description`** → to Filters to focus on top 5–10 locations

# Step 4: Format the Chart

- Click on **Color** → adjust palette to blue/orange like your chart
- Click on **Label** → check **"Show mark labels"** to display numbers
- Edit **Title**:
  **"Top Crime Locations in Chicago vs LA"**

### 5.4 Arrests per Crime Type

**Chicago**

```sql
SELECT
    Primary_Type,
    COUNT(*) AS total_crimes,
    SUM(CAST(Arrest AS INT)) AS arrests,
    ROUND(SUM(CAST(Arrest AS INT)) * 100.0 / COUNT(*), 2) AS
arrest_percentage
FROM crime_data_CH
GROUP BY Primary_Type
ORDER BY arrest_percentage DESC
LIMIT 5;
```

**Los Angeles**

```sql
SELECT
    Crm_Cd_Desc AS crime_type,
    COUNT(*) AS total_crimes,
    SUM(CASE WHEN Status IN ('AA', 'JA') THEN 1 ELSE 0 END) AS arrests,
    ROUND(SUM(CASE WHEN Status IN ('AA', 'JA') THEN 1 ELSE 0 END) * 100.0 /
COUNT(*), 2) AS arrest_percentage
FROM crime_data
GROUP BY Crm_Cd_Desc
ORDER BY arrest_percentage DESC
LIMIT 5;
```

## Visualize on Tableau for Arrests per Crime Type:

# STEP 1: Connect to the Dataset

1. Open **Tableau Desktop** or **Tableau Public**.
2. Connect to `crime_data_unified.csv`.
3. Go to **Sheet 1**.

## STEP 2: Build the Heat Map

**Drag fields into the worksheet:**

1. **Drag `crime_type`** → to **Rows**
2. **Drag `YEAR(occurrence_date)`** → to **Columns**
    - Right-click `occurrence_date` → Convert to **Discrete** → Choose **"Years"**
3. **Drag `city`** → to **Rows**, next to `crime_type`
4. **Drag `Number of Records`** → to **Color**
    - Use `SUM(Number of Records)`
5. **Drag `Number of Records`** → to **Label** so numbers appear inside each cell

## STEP 3: Format the Heat Map

1. On the **Marks card**, change type to **Square**
2. Adjust **Color**:
    - Click on **Color** → Choose a **green-yellow color scale** or **Orange-Green** diverging scale
3. Click **Label**:
    - Check **"Show mark labels"**
    - Format font size to ~9–10 for readability

## STEP 4: Filter and Focus

1. **Drag `city`** to **Filters** → select **Chicago** and **LA**
2. (Optional) **Drag `crime_type`** to Filters → show top 20 crime types only
    - Right-click `crime_type` → Filter → "Top" tab → "By field" → Top 20 by `SUM(Number of Records)`

## STEP 5: Add Chart Title

- Double-click title and type:
  **"Top Crime Types per Neighborhood Over Time"**

## Q2.Top Crime Types per Neighborhood Over Time

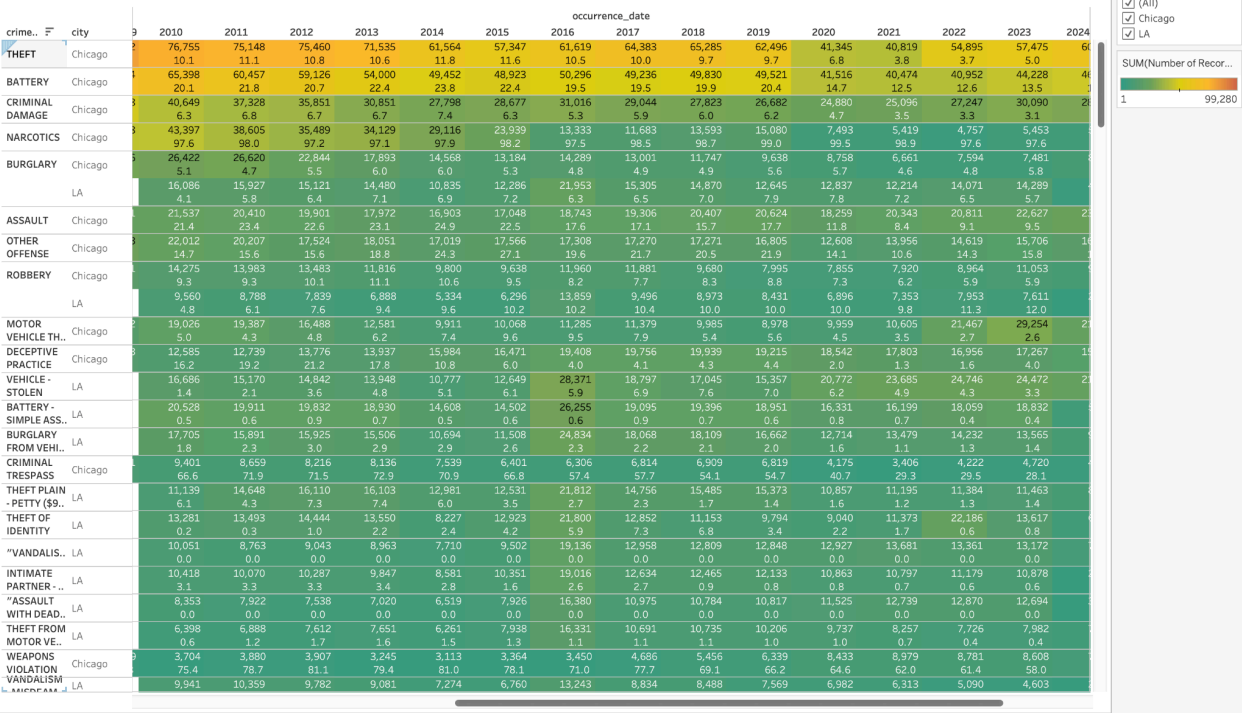| crime.. | city | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 | 2024 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| THEFT | Chicago | 76,755 | 75,148 | 75,460 | 71,535 | 61,564 | 57,347 | 61,619 | 64,383 | 65,285 | 62,496 | 41,345 | 40,819 | 54,895 | 57,475 | 6( |
| | | 10.1 | 11.1 | 10.8 | 10.6 | 11.8 | 11.6 | 10.5 | 10.0 | 9.7 | 9.7 | 6.8 | 3.8 | 3.7 | 5.0 | |
| BATTERY | Chicago | 65,398 | 60,457 | 59,126 | 54,000 | 49,452 | 48,923 | 50,296 | 49,236 | 49,830 | 49,521 | 41,516 | 40,474 | 40,952 | 44,228 | 4( |
| | | 20.1 | 21.8 | 20.7 | 22.4 | 23.8 | 22.4 | 19.5 | 19.5 | 19.9 | 20.4 | 14.7 | 12.5 | 12.6 | 13.5 | |
| CRIMINAL DAMAGE | Chicago | 40,649 | 37,328 | 35,851 | 30,851 | 27,798 | 28,677 | 31,016 | 29,044 | 27,823 | 26,682 | 24,880 | 25,096 | 27,247 | 30,090 | 2( |
| | | 6.3 | 6.8 | 6.7 | 6.7 | 7.4 | 6.3 | 5.3 | 5.9 | 6.0 | 6.2 | 4.7 | 3.5 | 3.3 | 3.1 | |
| NARCOTICS | Chicago | 43,397 | 38,605 | 35,489 | 34,129 | 29,116 | 23,939 | 13,333 | 11,683 | 13,593 | 15,080 | 7,493 | 5,419 | 4,757 | 5,453 | |
| | | 97.6 | 98.0 | 97.2 | 97.1 | 97.9 | 98.2 | 97.5 | 98.5 | 98.7 | 99.0 | 99.5 | 98.9 | 97.6 | 97.6 | |
| BURGLARY | Chicago | 26,422 | 26,620 | 22,844 | 17,893 | 14,568 | 13,184 | 14,289 | 13,001 | 11,747 | 9,638 | 8,758 | 6,661 | 7,594 | 7,481 | |
| | | 5.1 | 4.7 | 5.5 | 6.0 | 6.0 | 5.3 | 4.8 | 4.9 | 4.9 | 5.6 | 5.7 | 4.6 | 4.8 | 5.8 | |
| | LA | 16,086 | 15,927 | 15,121 | 14,480 | 10,835 | 12,286 | 21,953 | 15,305 | 14,870 | 12,645 | 12,837 | 12,214 | 14,071 | 14,289 | |
| | | 4.1 | 5.8 | 6.4 | 7.1 | 6.9 | 7.2 | 6.3 | 6.5 | 7.0 | 7.9 | 7.8 | 7.2 | 6.5 | 5.7 | |
| ASSAULT | Chicago | 21,537 | 20,410 | 19,901 | 17,972 | 16,903 | 17,048 | 18,743 | 19,306 | 20,407 | 20,624 | 18,259 | 20,343 | 20,811 | 22,627 | 2: |
| | | 21.4 | 23.4 | 22.6 | 23.1 | 24.9 | 22.5 | 17.6 | 17.1 | 15.7 | 17.7 | 11.8 | 8.4 | 9.1 | 9.5 | |
| OTHER OFFENSE | Chicago | 22,012 | 20,207 | 17,524 | 18,051 | 17,019 | 17,566 | 17,308 | 17,270 | 17,271 | 16,805 | 12,608 | 13,956 | 14,619 | 15,706 | 1( |
| | | 14.7 | 15.6 | 15.6 | 18.8 | 24.3 | 27.1 | 19.6 | 21.7 | 20.5 | 21.9 | 14.1 | 10.6 | 14.3 | 15.8 | |
| ROBBERY | Chicago | 14,275 | 13,983 | 13,483 | 11,816 | 9,800 | 9,638 | 11,960 | 11,881 | 9,680 | 7,995 | 7,855 | 7,920 | 8,964 | 11,053 | |
| | | 9.3 | 9.3 | 10.1 | 11.1 | 10.6 | 9.5 | 8.2 | 7.7 | 8.3 | 8.8 | 7.3 | 6.2 | 5.9 | 5.9 | |
| | LA | 9,560 | 8,788 | 7,839 | 6,888 | 5,334 | 6,296 | 13,859 | 9,496 | 8,973 | 8,431 | 6,896 | 7,353 | 7,953 | 7,611 | |
| | | 4.8 | 6.1 | 7.6 | 9.4 | 9.6 | 10.2 | 10.2 | 10.4 | 10.0 | 10.0 | 10.0 | 9.8 | 11.3 | 12.0 | |
| MOTOR VEHICLE TH.. | Chicago | 19,026 | 19,387 | 16,488 | 12,581 | 9,911 | 10,068 | 11,285 | 11,379 | 9,985 | 8,978 | 9,959 | 10,605 | 21,467 | 29,254 | 2: |
| | | 5.0 | 4.3 | 4.8 | 6.2 | 7.4 | 9.6 | 9.5 | 7.9 | 5.4 | 5.6 | 4.5 | 3.5 | 2.7 | 2.6 | |
| DECEPTIVE PRACTICE | Chicago | 12,585 | 12,739 | 13,776 | 13,937 | 15,984 | 16,471 | 19,408 | 19,756 | 19,939 | 19,215 | 18,542 | 17,803 | 16,956 | 17,267 | 1! |
| | | 16.2 | 19.2 | 21.2 | 17.8 | 10.8 | 6.0 | 4.0 | 4.1 | 4.3 | 4.4 | 2.0 | 1.3 | 1.6 | 4.0 | |
| VEHICLE - STOLEN | LA | 16,686 | 15,170 | 14,842 | 13,948 | 10,777 | 12,649 | 28,371 | 18,797 | 17,045 | 15,357 | 20,772 | 23,685 | 24,746 | 24,472 | 2: |
| | | 1.4 | 2.1 | 3.6 | 4.8 | 5.1 | 6.1 | 5.9 | 6.9 | 7.6 | 7.0 | 6.2 | 4.9 | 4.3 | 3.3 | |
| BATTERY - SIMPLE ASS.. | LA | 20,528 | 19,911 | 19,832 | 18,930 | 14,608 | 14,502 | 26,255 | 19,095 | 19,396 | 18,951 | 16,331 | 16,199 | 18,059 | 18,832 | |
| | | 0.5 | 0.6 | 0.9 | 0.7 | 0.5 | 0.6 | 0.6 | 0.9 | 0.7 | 0.6 | 0.8 | 0.7 | 0.4 | 0.4 | |
| BURGLARY FROM VEHI.. | LA | 17,705 | 15,891 | 15,925 | 15,506 | 10,694 | 11,508 | 24,834 | 18,068 | 18,109 | 16,662 | 12,714 | 13,479 | 14,232 | 13,565 | |
| | | 1.8 | 2.3 | 3.0 | 2.9 | 2.9 | 2.6 | 2.3 | 2.2 | 2.1 | 2.0 | 1.6 | 1.1 | 1.3 | 1.4 | |
| CRIMINAL TRESPASS | Chicago | 9,401 | 8,659 | 8,216 | 8,136 | 7,539 | 6,401 | 6,306 | 6,814 | 6,909 | 6,819 | 4,175 | 3,406 | 4,222 | 4,720 | |
| | | 66.6 | 71.9 | 71.5 | 72.9 | 70.9 | 66.8 | 57.4 | 57.7 | 54.1 | 54.7 | 40.7 | 29.3 | 29.5 | 28.1 | |
| THEFT PLAIN - PETTY ($9.. | LA | 11,139 | 14,648 | 16,110 | 16,103 | 12,981 | 12,531 | 21,812 | 14,756 | 15,485 | 15,373 | 10,857 | 11,195 | 11,384 | 11,463 | |
| | | 6.1 | 4.3 | 7.3 | 7.4 | 6.0 | 3.5 | 2.7 | 2.3 | 1.7 | 1.4 | 2.0 | 1.2 | 1.3 | 1.4 | |
| THEFT OF IDENTITY | LA | 13,281 | 13,493 | 14,444 | 13,550 | 8,227 | 12,923 | 21,800 | 12,852 | 11,153 | 9,794 | 9,040 | 11,373 | 22,186 | 13,617 | |
| | | 0.2 | 0.3 | 1.0 | 1.2 | 2.4 | 4.2 | 5.9 | 7.3 | 6.8 | 3.4 | 2.2 | 1.7 | 0.6 | 0.8 | |
| "VANDALIS.. | LA | 10,051 | 8,763 | 9,043 | 8,963 | 7,710 | 9,502 | 19,136 | 12,958 | 12,809 | 12,848 | 12,927 | 13,681 | 13,361 | 13,172 | |
| | | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | |
| INTIMATE PARTNER - .. | LA | 10,418 | 10,070 | 10,287 | 9,847 | 8,581 | 10,351 | 19,016 | 12,634 | 12,465 | 12,133 | 10,863 | 10,797 | 11,179 | 10,878 | |
| | | 3.1 | 3.3 | 3.3 | 3.4 | 2.8 | 1.6 | 2.6 | 2.7 | 0.9 | 0.8 | 0.7 | 0.6 | 0.6 | 0.6 | |
| "ASSAULT WITH DEAD.. | LA | 8,353 | 7,922 | 7,538 | 7,020 | 6,519 | 7,926 | 16,380 | 10,975 | 10,784 | 10,817 | 11,525 | 12,739 | 12,870 | 12,694 | |
| | | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | |
| THEFT FROM MOTOR VE.. | LA | 6,398 | 6,888 | 7,612 | 7,651 | 6,261 | 7,938 | 16,331 | 10,691 | 10,735 | 10,206 | 9,737 | 8,257 | 7,726 | 7,982 | |
| | | 0.6 | 1.2 | 1.7 | 1.6 | 1.5 | 1.3 | 1.1 | 1.1 | 1.1 | 1.0 | 0.7 | 0.4 | 0.4 | 0.4 | |
| WEAPONS VIOLATION | Chicago | 3,704 | 3,880 | 3,907 | 3,245 | 3,113 | 3,364 | 3,450 | 4,686 | 5,456 | 6,339 | 8,433 | 8,979 | 8,781 | 8,608 | |
| | | 75.4 | 78.7 | 81.1 | 79.4 | 81.0 | 78.1 | 71.0 | 77.7 | 69.1 | 66.2 | 64.6 | 62.0 | 61.4 | 58.0 | |
| VANDALISM MISDEM.. | LA | 9,941 | 10,359 | 9,782 | 9,081 | 7,274 | 6,760 | 13,243 | 8,834 | 8,488 | 7,569 | 6,982 | 6,313 | 5,090 | 4,603 | |

city
☑ (All)
☑ Chicago
☑ LA

SUM(Number of Recor...
1     99,280

# 5.5 Crimes Per Year

**Chicago**

```sql
SELECT
    Year,
    COUNT(*) AS crime_count
FROM crime_data_CH
WHERE Year IS NOT NULL AND Year BETWEEN 1900 AND 2100
GROUP BY Year
ORDER BY Year
LIMIT 5;
```

**Los Angeles**

```sql
SELECT
    SUBSTR(DATE_OCC, 7, 4) AS year, -- Extracts the 4-digit year (e.g.,
"2023") from DATE_OCC string, starting at position 7 (after "MM/DD/") for 4
characters

    COUNT(*) AS crime_count
FROM crime_data
WHERE DATE_OCC IS NOT NULL
GROUP BY SUBSTR(DATE_OCC, 7, 4)
ORDER BY year
LIMIT 5;
```

## Visualize on Tableau for Crimes Per Year:

# STEP 1: Load the Data

1. Open **Tableau Desktop** or **Tableau Public**.
2. Connect to your dataset: `crime_data_unified.csv`
3. Go to **Sheet 1**

# STEP 2: Drag Fields to the View

1. **Drag `YEAR(occurrence_date)`** → to **Columns**
   - Right-click `occurrence_date` → select **"Year"**
2. **Drag `Number of Records`** → to **Rows**
   - Use `SUM(Number of Records)`
3. **Drag `city`** → to **Color** on the **Marks** card
4. Optional: Drag `city` → to **Label** (to show name on line ends)
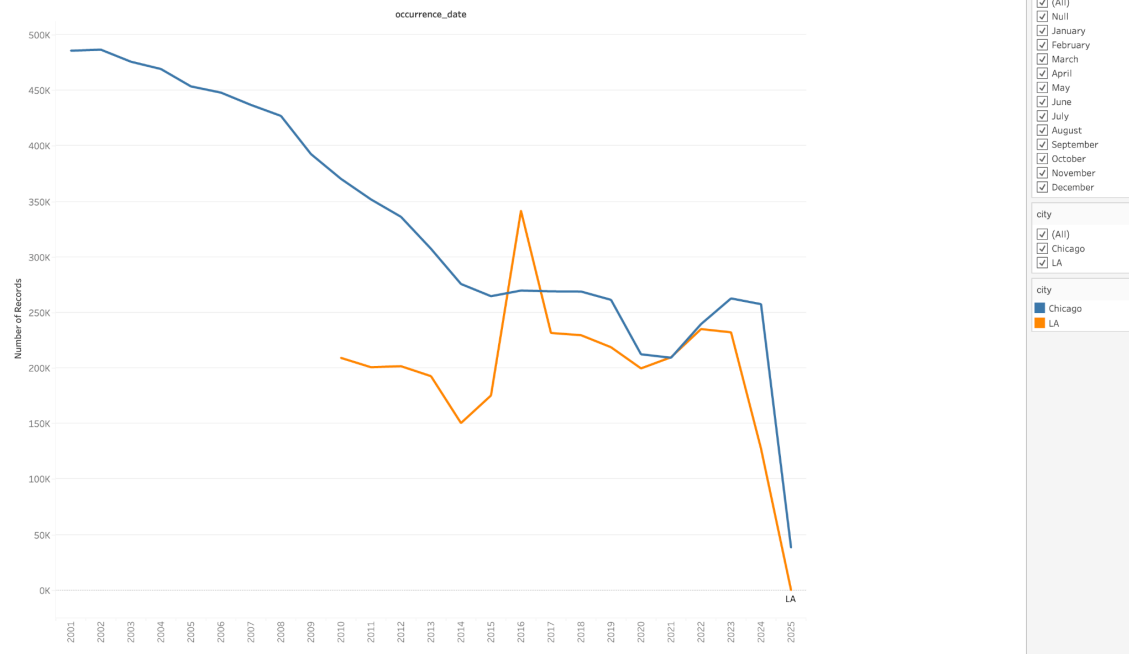
# STEP 3: Filter the Data

1. **Drag `city`** → to Filters → select **Chicago and LA**
2. Optional: Drag `MONTH(occurrence_date)` → to Filters (as you did) if you want to filter by months
3. Right-click filters and select **"Show Filter"** to add interactivity

# STEP 4: Format the Chart

1. Set chart type to **Line** (Tableau does this automatically)
2. Click **Label**:
    - Enable **Show mark labels** if you want point values
3. Add a title:
    **"Number of Crimes Over Time (Yearly) – Chicago vs LA"**

# References:

1. **URL of Data source:**

   **LA Crime Data (2020–Present)**:
   https://data.lacity.org/Public-Safety/Crime-Data-from-2020-to-Present/2nrs-mtv8/about_data
   **LA Crime Data (2010–2019)**:
   https://data.lacity.org/Public-Safety/Crime-Data-from-2010-to-2019/63jg-8b9z/about_data
   **Chicago Crime Data (2001–Present)**:
   https://data.cityofchicago.org/Public-Safety/Crimes-2025/t7ek-mgzi/about_data

2. **URL of GitHub: https://github.com/shjepz/BigDataCrimeAnalysis**