



CIS5200 Term Project Tutorial



Authors: Serena Tang, Sophathriya Sen, Alicia Martinez, Dima Kim

Instructor: Jongwook Woo

Date: 12/14/2024

Lab Tutorial

Serena Tang (xtang13@calstatela.edu), Sophathriya Sen (ssen5@calstatela.edu), Alicia Martinez (amarti761@calstatela.edu), Dima Kim (dkim171@calstatela.edu)

Crime Analysis - Chicago & Los Angeles

Objectives

- Get data from government open source
- Clean data using Python
- Create Hadoop cluster
- Visualization using Tableau and ArcGIS

Platform Specification

Cluster Version	Oracle Big Data Service
Number of Nodes	5
CPU Speed	2.5 GHz
Memory Size	30 GB
HDFS Capacity	1 TB
Hadoop Version	3.1.2
Hive Version	3.1.2

Step 1: Download and prepare Los Angeles and Chicago map data for visualization in ARCGIS

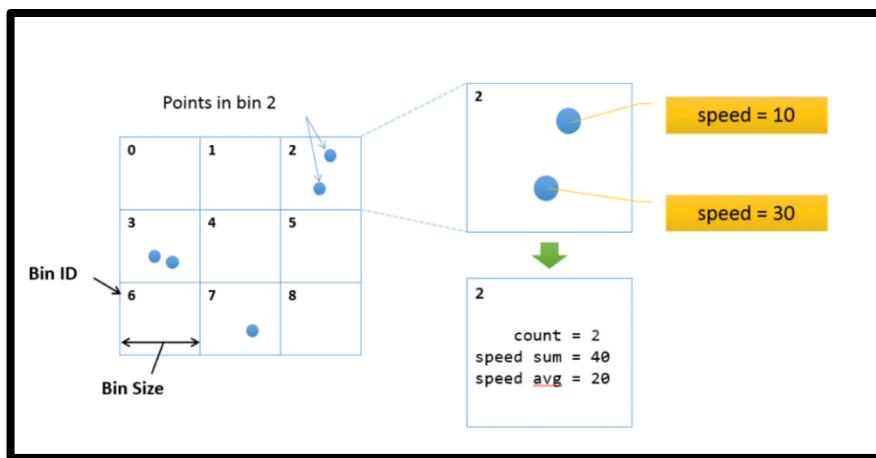
Download Los Angeles crime data set from 2020 to present from:

<https://catalog.data.gov/dataset/crime-data-from-2020-to-present>

Download Chicago crime data set from 2001 to present from:

<https://catalog.data.gov/dataset/crimes-2001-to-present>

Prepare Los Angeles and Chicago Crime data for spatial aggregation (sometimes called spatial binning). Spatial aggregation is extremely useful in summarizing big data to gain a meaningful snapshot of patterns in your data. Spatial aggregation works by creating square bins of a user specified size, like this:



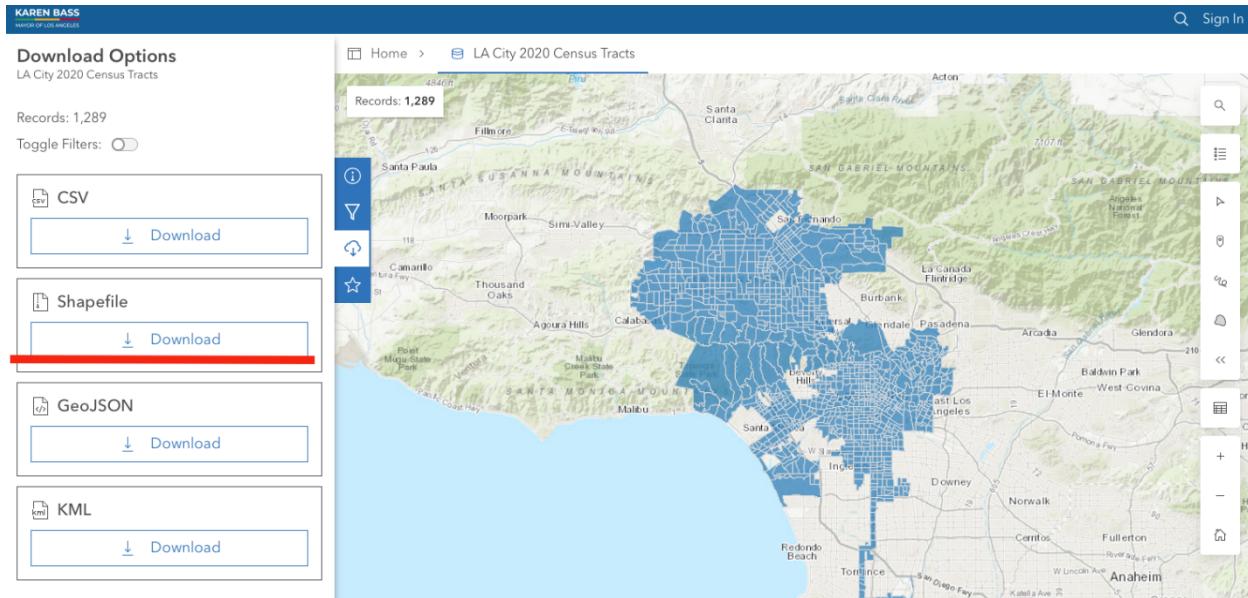
For the aggregation we need to download shapefiles that contain square bins for Los Angeles and Chicago.

Chicago:

<https://data.cityofchicago.org/api/geospatial/5jrd-6zik?method=export&format=Shapefile>

Los Angeles:

<https://geohub.lacity.org/datasets/la-city-2020-census-tracts-/explore?location=34.018934%2C-118.412043%2C10.11>

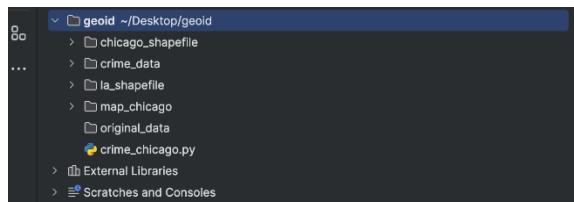


Next we are going to use Python to convert latitude and longitude from the original crime data files to square bins and create .csv files with four columns:

- geoid
- primary type
- Count
- Year

Chicago:

Open PyCharm IDE, then open folder where we saved all data and create the corresponding directories:



- Create a new Python file **crime_chicago.py**, paste the code below and execute it:

```

import pandas as pd
import geopandas as gpd
import zipfile

census_tract_shp = 'chicago_shapefile/geo_export_5c337cc7-e5d1-4343-a2ec-2aae0af46f7e.shp'
crime_csv = 'original_data/crime_chicago_2001_2024.csv'
# Generate list of years from 2020 to 2024

```

```

years = list(range(2020, 2025))
# Read in the census tract shapefile
tracts_gdf = gpd.read_file(census_tract_shp)
tracts_gdf = tracts_gdf.to_crs(epsg=4326)
tracts_gdf = tracts_gdf[['geometry', 'geoid10']]

# Read in the crime points data, this is very large dataset
crime_df = pd.read_csv(crime_csv)

# Show first 5 rows of tracts_gdf, you notice the geoid10 is 11 digits
tracts_gdf.head()

# Define function to aggregate crime points counts by census tract, by type of crime
def aggregate_points(points_gdf, geometry_gdf):
    # Spatial join the crime points with the census tracts
    points_gdf = gpd.sjoin(points_gdf, geometry_gdf, how='inner', predicate='within')
    # Group by census tract and type of crime
    points_gdf = points_gdf.groupby(['geoid10', 'Primary Type']).size().reset_index(name='count')
    return points_gdf

# Create an empty dataframe to store the combined results
combined_results = pd.DataFrame()

for year in years:
    # Filter the crime data for the current year
    crime_df_year = crime_df[crime_df['Year'] == year]

    # Convert the filtered crime data to a GeoDataFrame
    crime_gdf_year = gpd.GeoDataFrame(crime_df_year,
                                       geometry=gpd.points_from_xy(crime_df_year.Longitude, crime_df_year.Latitude))
    crime_gdf_year = crime_gdf_year.set_crs(epsg=4326)

    # Aggregate the crime points by census tract
    aggregated_points = aggregate_points(crime_gdf_year, tracts_gdf)

    # Add the year to the aggregated points dataframe
    aggregated_points['Year'] = year

    # Append the results to the combined dataframe
    combined_results = pd.concat([combined_results, aggregated_points], ignore_index=True)

# Save the combined results to a CSV file
combined_results.to_csv('crime_data/crime_chicago_2001_2024_by_tract_type.csv', index=False)
# Randomly sample 10,000 rows from the combined results, and display the first 5 rows.
combined_results.sample(10000).head()

# Filter combined_results where year is from 2020 to 2024
combined_results_2020_2024 = combined_results[combined_results['Year'] >= 2020]
# Summarize data by year and geoid10
combined_results_2020_2024 = combined_results_2020_2024.groupby(['geoid10', 'Year']).sum().reset_index()
# Drop primary type column
combined_results_2020_2024 = combined_results_2020_2024.drop(columns='Primary Type')
# Copy 'geometry' column from tracts_gdf to combined_results_2020_2024

```

```

combined_results_2020_2024 = combined_results_2020_2024.merge(tracts_gdf, on='geoid10', how='left')
# Convert combined_results_2020_2024 to GeoDataFrame
combined_results_2020_2024 = gpd.GeoDataFrame(combined_results_2020_2024, geometry='geometry')
# Change count, Year columns to integer
combined_results_2020_2024['count'] = combined_results_2020_2024['count'].astype(int)
# Save the combined results to a shapefile
combined_results_2020_2024.to_file('map_chicago/crime_chicago_aggregated.shp')

# Zip the shapefile
with zipfile.ZipFile('map_chicago/crime_chicago_aggregated.zip', 'w') as z:
    z.write('map_chicago/crime_chicago_aggregated.shp')
    z.write('map_chicago/crime_chicago_aggregated.shx')
    z.write('map_chicago/crime_chicago_aggregated.dbf')
    z.write('map_chicago/crime_chicago_aggregated.prj')
    z.write('map_chicago/crime_chicago_aggregated.cpg')

```

- Then check the new created files in the corresponding folders:

```

combined_results = pd.concat([combined_results, aggregated_points], ignore_index=True)
# Save the combined results to a CSV file
combined_results.to_csv(path_or_buf='crime_data/crime_chicago_2020_2024_by_tract_type.csv', index=False)
# Randomly sample 10,000 rows from the combined results, and display the first 5 rows
combined_results.sample(10000).head()

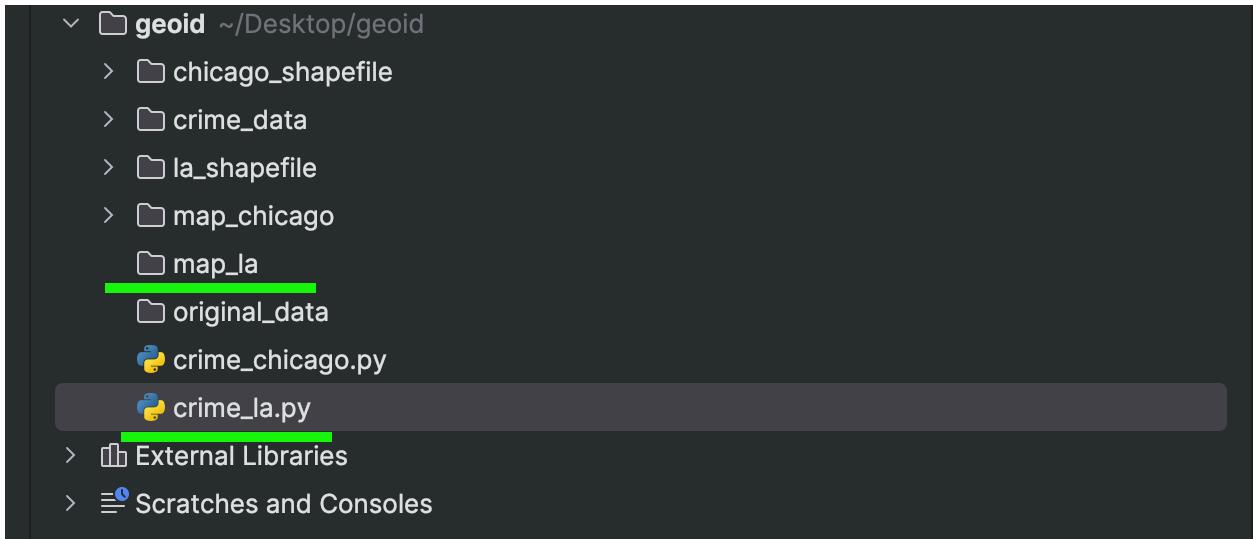
# Filter combined_results where year is from 2020 to 2024
combined_results_2020_2024 = combined_results[combined_results['Year'] >= 2020]
# Summarize data by year and geoid10
combined_results_2020_2024 = combined_results_2020_2024.groupby(['geoid10', 'Year']).sum().reset_index()
# Drop primary type column
combined_results_2020_2024 = combined_results_2020_2024.drop(columns='Primary Type')
# Copy geometry column from tracts_gdf to combined_results_2020_2024
combined_results_2020_2024 = combined_results_2020_2024.merge(tracts_gdf, on='geoid10', how='left')
# Convert combined_results_2020_2024 to GeoDataFrame
combined_results_2020_2024 = gpd.GeoDataFrame(combined_results_2020_2024, geometry='geometry')
# Change count, Year columns to integer
combined_results_2020_2024['count'] = combined_results_2020_2024['count'].astype(int)
# Save the combined results to a shapefile
combined_results_2020_2024.to_file('map_chicago/crime_chicago_aggregated.shp')

# Zip the shapefile
with zipfile.ZipFile(file='map_chicago/crime_chicago_aggregated.zip', mode='w') as z:
    z.write('map_chicago/crime_chicago_aggregated.shp')
    z.write('map_chicago/crime_chicago_aggregated.shx')
    z.write('map_chicago/crime_chicago_aggregated.dbf')
    z.write('map_chicago/crime_chicago_aggregated.prj')
    z.write('map_chicago/crime_chicago_aggregated.cpg')

```

Los Angeles:

- Create new folder map_la and then create a new Python file crime_la.py.



- Paste the code below and execute it:

```
import pandas as pd
import geopandas as gpd
import zipfile

census_tract_shp = 'la_shapefile/LA_City_2020_Census_Tracts.shp'
crime_csv = 'original_data/crime_la_2020_2024.csv'

# Generate list of years from 2001 to 2024
years = list(range(2020, 2024))
# Read in the census tract shapefile
tracts_gdf = gpd.read_file(census_tract_shp)
tracts_gdf = tracts_gdf.to_crs(epsg=4326)
tracts_gdf = tracts_gdf[['geometry', 'CT20']]

# Read in the crime points data, this is very large dataset
crime_df = pd.read_csv(crime_csv)

# Convert the date to a datetime object
crime_df['Date Rptd'] = pd.to_datetime(crime_df['Date Rptd'], format='%m/%d/%Y %I:%M:%S %p')
crime_df['Year'] = crime_df['Date Rptd'].dt.year

# Define function to aggregate crime points counts by census tract, by type of crime
def aggregate_points(points_gdf, geometry_gdf):
    # Spatial join the points to the census tracts
    points_gdf = gpd.sjoin(points_gdf, geometry_gdf, how='inner', predicate='within')
    # Group by the census tract and
    points_gdf = points_gdf.groupby(['CT20', 'Crm Cd']).size().reset_index(name='count')
    return points_gdf

# Create an empty dataframe to store the combined results
combined_results = pd.DataFrame()

# Loop through each year
```

```

for year in years:
    # Filter the crime data for the current year
    crime_df_year = crime_df[crime_df['Year'] == year]

    # Convert the filtered crime data to a GeoDataFrame
    crime_gdf_year = gpd.GeoDataFrame(crime_df_year, geometry=gpd.points_from_xy(crime_df_year.LON,
    crime_df_year.LAT))
    crime_gdf_year = crime_gdf_year.set_crs(epsg=4326)

    # Aggregate the crime points by census tract
    aggregated_points = aggregate_points(crime_gdf_year, tracts_gdf)

    # Add the year to the aggregated points dataframe
    aggregated_points['Year'] = year

    # Append the results to the combined dataframe
    combined_results = pd.concat([combined_results, aggregated_points], ignore_index=True)

# Create a 'geoid10' columns. goeid10 = '06037' + CT20
combined_results['geoid10'] = '06037' + combined_results['CT20'].astype(str)
combined_results = combined_results[['geoid10', 'Crm Cd', 'count', 'Year']]

# Save the combined results to a CSV file
combined_results.to_csv('crime_data/crime_la_2020_2024_by_tract_type.csv', index=False)

# Randomly sample 10,000 rows from the combined results, and display the first 5 rows.
combined_results.sample(10000).head()

# Filter combined_results where year is from 2020 to 2024
combined_results_2020_2024 = combined_results[combined_results['Year'] >= 2020]
# Summarize data by year and geoid10
combined_results_2020_2024 = combined_results_2020_2024.groupby(['geoid10', 'Year']).sum().reset_index()
# Drop primary type column
combined_results_2020_2024 = combined_results_2020_2024.drop(columns='Crm Cd')
# Copy 'geometry' column from tracts_gdf to combined_results_2020_2024
tracts_gdf['geoid10'] = '06037' + tracts_gdf['CT20'].astype(str)
combined_results_2020_2024 = combined_results_2020_2024.merge(tracts_gdf, on='geoid10', how='left')
# Convert combined_results_2020_2024 to GeoDataFrame
combined_results_2020_2024 = gpd.GeoDataFrame(combined_results_2020_2024, geometry='geometry')
# Change count, Year columns to integer
combined_results_2020_2024['count'] = combined_results_2020_2024['count'].astype(int)
# Save the combined results to a shapefile
combined_results_2020_2024.to_file('map_la/crime_la_aggregated.shp')

# Zip the shapefile
with zipfile.ZipFile('map_la/crime_la_aggregated.zip', 'w') as z:
    z.write('map_la/crime_la_aggregated.shp')
    z.write('map_la/crime_la_aggregated.shx')
    z.write('map_la/crime_la_aggregated.dbf')
    z.write('map_la/crime_la_aggregated.prj')
    z.write('map_la/crime_la_aggregated.cpg')

```

- Then check the new created files in the corresponding folders:

The screenshot shows a Jupyter Notebook environment with the following details:

- Project:** The sidebar displays a project structure with files like `geoid`, `chicago_shapefile`, `crime_data`, `map_chicago`, and `map_la`. The file `crime_la.py` is currently selected.
- Code Editor:** The main area contains Python code for processing crime data. The code includes:
 - Reading CSV files for Chicago and LA crime data.
 - Merging the data into a single DataFrame.
 - Sampling 10,000 rows from the combined results.
 - Filtering results for the year 2020-2024.
 - Grouping by geoid10 and year to summarize data.
 - Dropping the primary type column.
 - Copying the geometry column from the tracts_gdf to the combined_results DataFrame.
 - Merging the tracts_gdf with the combined_results DataFrame.
 - Converting the combined_results DataFrame to a GeoDataFrame.
 - Writing the results to a shapefile named `map_la/crime_la_aggregated.shp`.
 - Zippering the shapefile into a ZIP file named `map_la/crime_la_aggregated.zip`, containing files: `map_la/crime_la_aggregated.shp`, `map_la/crime_la_aggregated.shx`, `map_la/crime_la_aggregated.dbf`, `map_la/crime_la_aggregated.prj`, and `map_la/crime_la_aggregated.cpp`.
- Run Cell:** The bottom left shows the command run in the terminal: `/usr/local/bin/python3.12 /Users/mityakim/Desktop/geoid/crime_la.py`. The output indicates the process finished with exit code 0.

Step 2: Crime Data Loaded into Hadoop Clusters

Then follow the steps:

- create a new folder on your PC
 - put all files in a one **Crime.zip**
 - crime_chicago_2001_2024_by_tract_type.csv (result file from the step 1)
 - crime_la_2020_2024_by_tract_type.csv (result file from the step 1)
 - crime_la_2020_2024.csv (original data set LA)
 - crime_chicago_2001_2024.csv (original data set Chicago)
 - open terminal and using your path and username enter command below and then enter your password:

```
scp "/Users/mityakim/Desktop/Crime.zip" dkim171@129.146.230.230:/home/dkim171
```

- enter to the Linux server:

ssh dkim171@129.146.230.230

- enter command:

```
unzip Crime.zip
```

- You may check out all directories exist, and its files are created:

```
-bash-4.2$ ls
```

```
Last login: Sat Oct 26 08:49:38 on ttys001
+ ~ ssh dkim171@129.146.230.230
dkim171@129.146.230.230's password:
Last login: Sat Oct 26 17:28:17 2024 from syn-076-168-149-137.res.spectrum.com
-bash-4.2$ ls
00000_0 Austin baseball_salaries_2003.txt Crime.zip dictionary.tsv log_result.out minion_tweets.tar.gz senti_out.csv Tweets
90972b.log Austin_Crime.csv Chicago Dallas LA __MACOSX SensorFiles time_zone_map.tsv
-bash-4.2$
```

- Create crime folders and sub folders

```
-bash-4.2$ hdfs dfs -mkdir Project/crime
-bash-4.2$ hdfs dfs -mkdir Project/crime/chicago
-bash-4.2$ hdfs dfs -mkdir Project/crime/la
-bash-4.2$ hdfs dfs -mkdir Project/crime/chicago_tractid
-bash-4.2$ hdfs dfs -mkdir Project/crime/la_tractid
```

- Place the files in Hadoop:

```
-bash-4.2$ hdfs dfs -put crime_chicago_2001_2024.csv Project/crime/chicago
-bash-4.2$ hdfs dfs -put crime_la_2020_2024.csv Project/crime/la
-bash-4.2$ hdfs dfs -put crime_chicago_2001_2024_by_tract_type.csv
Project/crime/chicago_tractid
-bash-4.2$ hdfs dfs -put crime_la_2020_2024_by_tract_type.csv
Project/crime/la_tractid
```

- to check files, enter:

```
-bash-4.2$ hdfs dfs -ls Project/crime/chicago_tractid
-bash-4.2$ hdfs dfs -ls Project/crime/la_tractid
-bash-4.2$ hdfs dfs -ls Project/crime/chicago
-bash-4.2$ hdfs dfs -ls Project/crime/la
```

```
[bash-4.2$ hdfs dfs -ls Project/crime/chicago
Found 1 items
-rw-r--r-- 3 dkim171 hdfs 1932368003 2024-11-05 01:45 Project/crime/chicago/Crimes_-_2001_to_Present_20241019.csv
-bash-4.2$ ]
```

Step 3: Creating Hive tables to Query Crime Data

Open Hive CLI (Command Line Shell Interface):

```
-bash-4.2$ beeline
```

Make sure that you are using your database:

```
use <username>
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> use <username>
```

Los Angeles Crime Tractid Table

```
DROP TABLE IF EXISTS original_la_crime;

CREATE EXTERNAL TABLE IF NOT EXISTS original_la_crime(
    DR_NO STRING,
    Date_Rpt STRING,
    DATE_OCC STRING,
    TIME_OCC STRING,
    AREA STRING,
    AREA_NAME STRING,
    RptDist_No STRING,
    Part_1_2 STRING,
    Crm_Cd STRING,
    Crm_Cd_Desc STRING,
    Mocodes STRING,
    Vict_Age STRING,
    Vict_Sex STRING,
    Vict_Descent STRING,
    Premis_Cd STRING,
    Premis_Desc STRING,
    Weapon_Used_Cd STRING,
    Weapon_Desc STRING,
    Status STRING,
    Status_Desc STRING,
    Crm_Cd_1 STRING,
    Crm_Cd_2 STRING,
    Crm_Cd_3 STRING,
    Crm_Cd_4 STRING,
    LOCATION STRING,
    Cross_Street STRING,
    LAT STRING,
    LON STRING
)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
WITH SERDEPROPERTIES (
    "separatorChar" = ",",
    "quoteChar" = "\"",
    "escapeChar" = "\\"
)
LOCATION '/user/dkim171/Project/crime/la'
```

```
TBLPROPERTIES ('skip.header.line.count'='1');
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
+-----+  
| tab_name |  
+-----+  
| la_crime |  
+-----+
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from original_la_crime  
LIMIT 10;
```

Chicago Crime Table

```
DROP TABLE IF EXISTS original_chicago_crime;  
CREATE EXTERNAL TABLE IF NOT EXISTS original_chicago_crime(  
    ID STRING,  
    Case_Number STRING,  
    `Date` STRING,  
    Block STRING,  
    IUCR STRING,  
    Primary_Type STRING,  
    Description STRING,  
    Location_Description STRING,  
    Arrest STRING,  
    Domestic STRING,  
    Beat STRING,  
    District STRING,  
    Ward STRING,  
    Community_Area STRING,  
    FBI_Code STRING,  
    X_Coordinate STRING,  
    Y_Coordinate STRING,  
    Year STRING,  
    Updated_On STRING,  
    Latitude STRING,  
    Longitude STRING,  
    Location STRING  
)  
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'  
WITH SERDEPROPERTIES (  
    "separatorChar" = ",",  
    "quoteChar" = "\"",  
    "escapeChar" = "\\\"
```

```
)  
LOCATION '/user/dkim171/Project/crime/chicago'  
TBLPROPERTIES ('skip.header.line.count'='1');
```

- Check if the table was successfully created:
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> **show tables;**

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from original_la_crime  
LIMIT 10;
```

Los Angeles Crime Tractid Table

```
DROP TABLE IF EXISTS la_tractid;  
  
CREATE EXTERNAL TABLE IF NOT EXISTS la_tractid(  
    tract_id STRING,  
    primary_type STRING,  
    count STRING,  
    year STRING  
)  
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'  
WITH SERDEPROPERTIES (  
    "separatorChar" = ",",  
    "quoteChar" = "\"",  
    "escapeChar" = "\\\""  
)  
LOCATION '/user/dkim171/Project/crime/la_tractid'  
TBLPROPERTIES ('skip.header.line.count'='1');
```

- Check if the table was successfully created:
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> **show tables;**

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from la_tractid LIMIT 10;
```

Chicago Crime Tractid Table

```
DROP TABLE IF EXISTS chicago_tractid;  
  
CREATE EXTERNAL TABLE IF NOT EXISTS chicago_tractid(  
    tract_id STRING,  
    primary_type STRING,  
    count STRING,  
    year STRING  
)  
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'  
WITH SERDEPROPERTIES (
```

```

"separatorChar" = ",",
"quoteChar" = "\"",
"escapeChar" = "\\"
)
LOCATION '/user/dkim171/Project/crime/chicago_tractid'
TBLPROPERTIES ('skip.header.line.count'=1');

```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from chicago_tractid
LIMIT 10;
```

Step 4: Download and clean topic data

Education Data

- Go to the website:

<https://data.census.gov/all>

For the analysis purposes we will use:

1. Geographies → Census Tract → California → Los Angeles County
→ All Census Tracts

The screenshot shows the data.census.gov website interface. At the top, there's a navigation bar with links for 'All', 'Tables', 'Maps', 'Profiles', and 'Pages'. Below the navigation is a search bar with a magnifying glass icon and a link to 'Advanced Search'. On the left, there's a sidebar titled '1 Filter' with a 'Filters' button. It shows a single filter applied: 'All Census Tracts within Los Angeles County, California'. There's also a 'Clear all filters' button and a search bar for filters. The main content area has a breadcrumb trail: 'California / Los Angeles County, California / Select Census Tract'. A dropdown menu titled 'Within other geographies' is open, showing a list of census tracts. The first item in the list, 'All Census Tracts within Los Angeles County, California', is highlighted with a blue border and checked. Other items in the list include: Census Tract 1011.10; Los Angeles County; California; Census Tract 1011.20, Los Angeles County, California; Census Tract 1011.22; Los Angeles County; California; Census Tract 1012.10, Los Angeles County, California; Census Tract 1012.20; Los Angeles County; California; Census Tract 1012.21; Los Angeles County; California; Census Tract 1012.22; Los Angeles County; California; Census Tract 1013; Los Angeles County; California; Census Tract 1014; Los Angeles County; California; Census Tract 1021.01, Los Angeles County, California; Census Tract 1021.02, Los Angeles County, California; and Census Tract 1021.03; Los Angeles County; California.

- Topics → Education → Educational Attainment

- Years → 2022

The screenshot shows the Census Bureau's data search interface. In the search bar, 'Educational Attainment' is entered. Below the search bar, the 'Tables' tab is selected. On the left, a sidebar displays filters for 'Geographies' (All Census Tracts within Los Angeles County, California) and 'Topics' (Education, Employment, Families and Living Arrangements, Health, Housing, Income and Poverty, Populations and People, Race and Ethnicity, Surveys, American Community Survey, Years: 2022, 2021, 2020). The main search results area shows 44 results for 'Educational Attainment' in 2022, including various tables from the American Community Survey (e.g., B1501, B1502, B1503, B1504, B1505, B1506).

Then we need to open Tables and click Download and repeat for 2020, 2021 years.

This screenshot shows the detailed view of the 'B1501 | Educational Attainment' table. The table is described as 'American Community Survey | 2022 ACS 5-Year Estimates Subject Tables'. It includes columns for 'View: 10 25 50' and 'Download Table Data'. A large blue button labeled 'Download Table Data' is prominent. To the right, there is a message about the table being too large for the browser and a link to 'Technical Details'. At the bottom, there is a link to 'open the table'.

- After we need to clean our data and get columns that we will use for the future analysis, and we are going to use Python and PyCharm or Visual Studio Code.

```
import pandas as pd
```

```
# Define the mapping of original columns to new column names
```

```

column_map = {
    "GEO_ID": "geo_id",
    "NAME": "area_name",

    "S1501_C01_001E": "population_age_18_24",
    "S1501_C01_003E": "high_school_grad_or_higher_18_24",

    "S1501_C01_016E": "population_25_34_years",
    "S1501_C01_017E": "high_school_grad_or_higher_25_34",

    "S1501_C01_019E": "population_35_44_years",
    "S1501_C01_020E": "high_school_grad_or_higher_35_44",

    "S1501_C01_022E": "population_45_64_years",
    "S1501_C01_023E": "high_school_grad_or_higher_45_64",

    "S1501_C01_025E": "population_age_65_and_over",
    "S1501_C01_026E": "high_school_grad_or_higher_65_and_over"
}

# Load the data from the provided file
file_path = '/Users/mityakim/Desktop/education/2020_education.csv' # replace with the actual
file path
data = pd.read_csv(file_path)

# Select only the required columns and rename them
selected_data = data[list(column_map.keys())].copy()
selected_data.rename(columns=column_map, inplace=True)

# Add the new 'year' column with a constant value of 2020
selected_data['year'] = 2020

# Save the cleaned and selected data to a new CSV file
output_file_path = '/Users/mityakim/Desktop/education/2020_education_clean.csv' # replace
with desired output path
selected_data.to_csv(output_file_path, index=False)

print("CSV file created successfully at:", output_file_path)

```

3. Repeat operation for all 3 files (2020, 2021, 2022).

```

Project
python ~/Desktop/education/python
  main.py
...
External Libraries
Scratches and Consoles

main.py
26  data = pd.read_csv(file_path)
27
28  # Select only the required columns and rename them
29  selected_data = data[list(column_map.keys())].copy()
30  selected_data.rename(columns=column_map, inplace=True)
31
32  # Add the new 'year' column with a constant value of 2021
33  selected_data['year'] = 2022
34
35  # Save the cleaned and selected data to a new CSV file
36  output_file_path = '/Users/mityakim/Desktop/education/2022_education_clean.csv' # replace with desired output path
37  selected_data.to_csv(output_file_path, index=False)
38
39  print("CSV file created successfully at:", output_file_path)
40

```

Run main

```

/usr/local/bin/python3.12 /Users/mityakim/Desktop/education/python/main.py
[2]: D:\Users\mityakim\Desktop\education\main.py:24: DtypeWarning: Columns (2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37,38,39,40,41,42,43,
      data = pd.read_csv(file_path)
CSV file created successfully at: /Users/mityakim/Desktop/education/2022_education_clean.csv

Process finished with exit code 0

```

python main.py

33:29 LF UTF-8 4 spaces Python 3.12

- Then change filter to Chicago

United States Census Bureau

Tables

3 Filters

All

Geographies

Topics

Years

Illinois / Cook County, Illinois / Select Census Tract

Within other geographies

Search Census Tract

All Census Tracts within Cook County, Illinois

Census Tract 0, Cook County, Illinois

Census Tract 101, Cook County, Illinois

Census Tract 102, Cook County, Illinois

Census Tract 102.01, Cook County, Illinois

Census Tract 102.02, Cook County, Illinois

Census Tract 103, Cook County, Illinois

Census Tract 104, Cook County, Illinois

Census Tract 105, Cook County, Illinois

Census Tract 105.01, Cook County, Illinois

Census Tract 105.02, Cook County, Illinois

Census Tract 105.03, Cook County, Illinois

Census Tract 106, Cook County, Illinois

Census Tract 107, Cook County, Illinois

Census Tract 107.01, Cook County, Illinois

Census Tract 107.02, Cook County, Illinois

Census Tract 108, Cook County, Illinois

Census Tract 109, Cook County, Illinois

Census Tract 201, Cook County, Illinois

Census Tract 202, Cook County, Illinois

Census Tract 203, Cook County, Illinois

Census Tract 203.01, Cook County, Illinois

Census Tract 203.02, Cook County, Illinois

Census Tract 204, Cook County, Illinois

Census Tract 205, Cook County, Illinois

Census Tract 206, Cook County, Illinois

Census Tract 206.01, Cook County, Illinois

44 Results

View: 10 | 25 | 50 Download Table Data

American Community Survey
B15001 | Educational Attainment
2022: ACS 5-Year Estimates Subject: Tables

American Community Survey
B1502 | Field of Bachelor's Degree for First Major
2022: ACS 5-Year Estimates Subject: Tables

American Community Survey
B06009 | Place of Birth by Educational Attainment in the Past Year
2022: ACS 5-Year Estimates Detailed Tables

American Community Survey
B07009 | Geographical Mobility in the Past Year by Educational Attainment
2022: ACS 5-Year Estimates Detailed Tables

American Community Survey
B13014 | Women 15 to 50 Years Who Had a Birth in the Past Year by Educational Attainment
2022: ACS 5-Year Estimates Detailed Tables

American Community Survey
B15001 | Sex by Educational Attainment for the Population
2022: ACS 5-Year Estimates Detailed Tables

American Community Survey
B15002 | Sex by Educational Attainment for the Population
2022: ACS 5-Year Estimates Detailed Tables

American Community Survey
B15003 | Educational Attainment for the Population 25 Years and Older
2022: ACS 5-Year Estimates Detailed Tables

American Community Survey
B15011 | Sex by Age by Field of Bachelor's Degree
2022: ACS 5-Year Estimates Detailed Tables

Download 2020, 2021, 2022 years and repeat steps for 2 and 3.

Employment Data

- Topics → Employment → Employment and Labor Force Status

The screenshot shows the Census Bureau's website interface. In the top left, there's a note: "An official website of the United States government. How's how you know". The top navigation bar includes "Search", "Advanced Search", and tabs for "All", "Tables", "Maps", "Profiles", and "Pages". On the left, a sidebar titled "Filters" shows three selected filters: "2022", "All Census Tracts within Cook County, Illinois", and "Employment and Labor Force Status". Below this is a search bar and a "Clear all filters" button. The main content area is titled "Select Employment" and lists several categories: "Work Experience" (unchecked), "Class of Worker" (unchecked), "Commuting" (unchecked), "Employment" (unchecked), "Employment and Labor Force Status" (checked), "Industry" (unchecked), and "Occupation" (unchecked). To the right, a large panel displays "100 Results" for the "Employment and Labor Force Status" filter. It lists items such as "DP03 | Selected Economic Characteristics", "S2301 | Employment Status", "S2302 | Employment Characteristics of Families", "S2303 | Work Status in the Past 12 Months", "S2401 | Occupation by Sex for the Civilian Employed Full-Time, Year-Round", "S2402 | Occupation by Sex for the Full-Time, Year-Round", and "S2403 | Industry for the Civilian Employment". A "Download Table Data" button is visible at the top of this panel.

The same filters for the Year and Geographic. Download data for LA and Chicago for 2020, 2021 and 2022.

- Next clean data and get columns using Python:

```
import pandas as pd

# Define the mapping of original columns to new column names
column_map = {
    "GEO_ID": "geo_id",
    "NAME": "area_name",

    "DP03_0001E": "population_age_16_and_over",
    "DP03_0002E": "population_age_16_and_over_in_labor_force",
    "DP03_0004E": "population_age_16_and_over_in_civ_labor_force_employed",
    "DP03_0005E": "population_age_16_and_over_in_civ_labor_force_unemployed",
    "DP03_0007E": "population_age_16_and_over_not_in_labor_force",
    "DP03_0014E": "own_children_of_the_householder_under_6_years",
    "DP03_0016E": "own_children_of_the_householder_under_6_to_17_years",
}

# Load the data from the provided file
file_path = '/Users/mityakim/Desktop/employment/2020_employment.csv' # replace with the actual file path
data = pd.read_csv(file_path)

# Select only the required columns and rename them
selected_data = data[list(column_map.keys())].copy()
selected_data.rename(columns=column_map, inplace=True)

# Add the new 'year' column with a constant value of 2020
```

```

selected_data['year'] = 2020

# Save the cleaned and selected data to a new CSV file
output_file_path = '/Users/mityakim/Desktop/employment/2020_employment_clean.csv' # replace with desired output path
selected_data.to_csv(output_file_path, index=False)

print("CSV file created successfully at:", output_file_path)

```

The screenshot shows a Python development environment with a project structure on the left and a code editor on the right. The code editor contains a file named 'main.py' with the following content:

```

14     "DP03_0016M": "own_children_of_the_householder_under_6_to_17_years",
15 }
16
17 # Load the data from the provided file
18 file_path = '/Users/mityakim/Desktop/employment/2020_employment.csv' # replace with the actual file path
19 data = pd.read_csv(file_path)
20
21 # Select only the required columns and rename them
22 selected_data = data[list(column_map.keys())].copy()
23 selected_data.rename(columns=column_map, inplace=True)
24
25 # Add the new 'year' column with a constant value of 2021
26 selected_data['year'] = 2020
27
28 # Save the cleaned and selected data to a new CSV file
29 output_file_path = '/Users/mityakim/Desktop/employment/2020_employment_clean.csv' # replace with desired output path
30 selected_data.to_csv(output_file_path, index=False)
31
32 print("CSV file created successfully at:", output_file_path)
33

```

Below the code editor is the 'Run' tab, which is active. The run history shows the command run was '/usr/local/bin/python3.12 /Users/mityakim/Desktop/employment/python/main.py'. The output of the run is:

```

/usr/local/bin/python3.12 /Users/mityakim/Desktop/employment/python/main.py
/Users/mityakim/Desktop/employment/main.py:19: DtypeWarning: Columns (2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,37
  data = pd.read_csv(file_path)
CSV file created successfully at: /Users/mityakim/Desktop/employment/2020_employment_clean.csv

Process finished with exit code 0

```

The status bar at the bottom indicates the time is 26:29, the file encoding is UTF-8, there are 4 spaces, and the Python version is 3.12.

Income Data

- Topics → Income and Poverty → Income and Poverty

The same filters for the Year and Geographic. Download data for LA and Chicago for 2020, 2021 and 2022.

Next clean data and get columns using Python:

```
import pandas as pd

# Define the mapping of original columns to new column names
column_map = {
    "GEO_ID": "geo_id",
    "NAME": "area_name",

    "S1901_C01_002E": "less_than_10k",
    "S1901_C01_003E": "from_10k_to_14999",
    "S1901_C01_004E": "from_15k_to_24999",
    "S1901_C01_005E": "from_25k_to_34999",
    "S1901_C01_006E": "from_35k_to_49999",
    "S1901_C01_007E": "from_50k_to_74999",
    "S1901_C01_008E": "from_75k_to_99999",
    "S1901_C01_009E": "from_100k_to_149999",
    "S1901_C01_010E": "from_150k_to_199999",
    "S1901_C01_011E": "from_200k_or_more",
    "S1901_C01_012E": "median_income",
}

# Load the data from the provided file
file_path = '/Users/mityakim/Desktop/income/chicago/2022_income_chicago.csv' # replace with the actual file path
```

```

data = pd.read_csv(file_path)

# Select only the required columns and rename them
selected_data = data[list(column_map.keys())].copy()
selected_data.rename(columns=column_map, inplace=True)

# Add the new 'year' column with a constant value of 2021
selected_data['year'] = 2022

# Save the cleaned and selected data to a new CSV file
output_file_path =
'/Users/mityakim/Desktop/income/chicago/clean_chicago_income/2022_income_clean_chicago.csv' # replace with desired output path
selected_data.to_csv(output_file_path, index=False)

print("CSV file created successfully at:", output_file_path)

```

```

import pandas as pd

# Define the mapping of original columns to new column names
column_map = {
    "E001ID": "geo_id",
    "NAME": "area_name",
    ...
    "S1991_C01_002E": "less_than_10k",
    "S1991_C01_003E": "from_10K_to_14999",
    "S1991_C01_004E": "from_15K_to_24999",
    "S1991_C01_005E": "from_25K_to_34999",
    "S1991_C01_006E": "from_35K_to_49999",
    "S1991_C01_007E": "from_50K_to_74999",
    "S1991_C01_008E": "from_75K_to_104999",
    "S1991_C01_009E": "from_105K_to_149999",
    "S1991_C01_010E": "from_150K_to_199999",
    "S1991_C01_011E": "from_200K_on_more",
    "S1991_C01_012E": "median_income",
}

# Load the data from the provided file
file_path = '/Users/mityakim/Desktop/income/chicago/2022_income_chicago.csv' # replace with the actual file path
data = pd.read_csv(file_path)

# Select only the required columns and rename them
selected_data = data[list(column_map.keys())].copy()
selected_data.rename(columns=column_map, inplace=True)

# Add the new 'year' column with a constant value of 2021
selected_data['year'] = 2022

# Save the cleaned and selected data to a new CSV file
selected_data.to_csv(output_file_path, index=False)

print("CSV file created successfully at:", output_file_path)

```

Population Data

Topics → Population and People

- DP05 ACS Demographic and Housing Estimates

Label	Estimate	Percent	Estimate	Percent
SEX AND AGE				
Total population	4,014	4.014	4,164	4,164
Male	2,005	47.5%	2,098	50.4%
Female	2,009	52.5%	2,066	49.6%
Sex ratio (males per 100 females)	90.3	0.0	101.5	(X)
Under 5 years	141	3.5%	128	3.1%
5 to 9 years	154	3.8%	179	4.3%
10 to 14 years	209	5.0%	242	5.8%
15 to 19 years	241	6.0%	180	4.3%
20 to 24 years	179	4.5%	303	7.3%
25 to 34 years	325	13.2%	282	6.8%
35 to 44 years	533	13.3%	453	10.8%
45 to 54 years	633	15.4%	684	16.4%
55 to 59 years	256	6.4%	451	10.8%
60 to 64 years	322	8.0%	352	8.5%
65 to 74 years	603	15.0%	511	12.3%
75 to 84 years	203	5.1%	249	6.0%
85 years and over	34	0.8%	169	3.6%
Median age (years)	45.4	0.0	(X)	(X)
Under 10 years	1,372	37.7%	3,397	86.4%
10 years and over	3,419	62.2%	7,099	13.6%
21 years and over	3,298	82.2%	5,411	81.9%
62 years and over	1,049	26.4%	1,066	23.6%
65 years and over	840	20.9%	909	21.6%
EDUCATION				
18 years and over	3,119	3,119	3,499	3,499
Male	1,669	45.9%	1,790	51.2%
Female	1,850	54.1%	1,709	48.8%
Sex ratio (males per 100 females)	84.8	0.0	104.7	0.0
RACE AND ETHNICITY				
18 years and over	840	84.0%	808	80.8%
Male	389	46.3%	406	49.3%
Female	451	53.7%	503	50.7%

The same filters for the Year and Geographic. Download data for LA and Chicago for 2020, 2021 and 2022.

Next clean data and get columns using Python:

```
import pandas as pd

# Define the mapping of original columns to new column names
column_map = {
    "GEO_ID": "geo_id",
    "NAME": "area_name",

    "DP05_0036E": "one_race",
    "DP05_0038E": "black_or_African_american",
    "DP05_0039E": "american_indian_and_alaska_native",
    "DP05_0040E": "cherokee_tribal_grouping",
    "DP05_0041E": "chippewa_tribal_grouping",
    "DP05_0042E": "navajo_tribal_grouping",
    "DP05_0043E": "sioux_tribal_grouping",
    "DP05_0044E": "asian",
    "DP05_0045E": "indian",
    "DP05_0046E": "chinese",
    "DP05_0047E": "filipino",
    "DP05_0048E": "japanese",
    "DP05_0049E": "korean",
    "DP05_0050E": "vietnamese",
    "DP05_0051E": "other_asian",
```

```

"DP05_0052E": "native_hawaiian_and_other_pacific_islander",
"DP05_0053E": "chamorro",
"DP05_0054E": "native_hawaiian",
"DP05_0055E": "samoan",
"DP05_0056E": "other_native_hawaiian_and_other_pacific_islander",
"DP05_0065E":
"race_alone_or_in_combination_with_one_or_more_other_races_total_population",
"DP05_0072E": "hispanic_or_latino_total_population",
"DP05_0073E": "hispanic_or_latino_total_population_of_any_race",
"DP05_0074E": "mexican",
"DP05_0075E": "puerto_rican",
"DP05_0076E": "cuban",
"DP05_0076M": "other_hispanic_or_latino"
}

# Load the data from the provided file
file_path = '/Users/mityakim/Desktop/population/la/2020_population_la.csv'
data = pd.read_csv(file_path)

# Select only the required columns and rename them
selected_data = data[list(column_map.keys())].copy()
selected_data.rename(columns=column_map, inplace=True)

# Add the new 'year' column with a constant value of 2020
selected_data['year'] = 2020

# Save the cleaned and selected data to a new CSV file
output_file_path =
'/Users/mityakim/Desktop/population/la/clean_la_population/2020_population_clean_la.csv'
selected_data.to_csv(output_file_path, index=False)

print("CSV file created successfully at:", output_file_path)

```

```

import pandas as pd

# Define the mapping of original columns to new column names
column_map = {
    "GEO_ID": "geo_id",
    "NAME": "area_name",

    "DP05_0036E": "one_race",
    "DP05_0038E": "black_or_African_american",
    "DP05_0039E": "american_indian_and_alaska_native",
    "DP05_0040E": "cherokee_tribal_grouping",
    "DP05_0041E": "chippewa_tribal_grouping",
    "DP05_0042E": "navajo_tribal_grouping",
    "DP05_0043E": "sioux_tribal_grouping",
    "DP05_0044E": "asian",
    "DP05_0045E": "indian",
    "DP05_0046E": "hispanic"
}

```

Run main

```

/usr/local/bin/python3.12 /Users/mityakim/Desktop/population/main.py
/Users/mityakim/Desktop/population/main.py:39: DtypeWarning: Columns (2,3,4,5,6,7,10,11,12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34)
  data = pd.read_csv(file_path)
CSV file created successfully at: /Users/mityakim/Desktop/population/la/clean_la_population/2022_population_clean_la.csv
Process finished with exit code 0

```

Step 5: Education, Employment, Income, Population Data Loaded into Hadoop Clusters

- Put all clean files in data.zip and then open the terminal. Do not forget to use your path and username:

```
scp "/Users/mityakim/Desktop/data.zip" dkim171@129.146.230.230:/home/dkim171
```

```
[~] ~ scp "/Users/mityakim/Desktop/data.zip" dkim171@129.146.230.230:/home/dkim171
[dkim171@129.146.230.230's password:
data.zip                                              100%  569KB 322.8KB/s   00:01
[~] ~ ]
```

- Check files on Linux server:

```
ssh dkim171@129.146.230.230
```

- Enter command:

```
-bash-4.2$ ls
[-bash-4.2$ ls
00000_0           CrimesChicago.zip  log_result.out      SensorFiles
909f2b.log        Crime.zip       __MACOSX
Austin            Dallas          minion_tweets.tar.gz  senti_out.csv
Austin_Crime.csv  data.zip       dictionary.tsv    time_zone_map.tsv
baseball_salaries_2003.txt   LA             ratings_2012.txt  Tweets
Chicago          LOAD           reddit_opinion_PSE_ISR.csv
Crimes_-_2001_to_Present_20241019.csv  reddit.zip
-bash-4.2$ ]
```

- Then unzip data:

```
-bash-4.2$ unzip data
```

- Create project folder and sub folders

```
-bash-4.2$ hdfs dfs -mkdir Project
```

```
-bash-4.2$ hdfs dfs -mkdir Project/education
-bash-4.2$ hdfs dfs -mkdir Project/education/la
-bash-4.2$ hdfs dfs -mkdir Project/education/chicago
```

```
-bash-4.2$ hdfs dfs -mkdir Project/employment
-bash-4.2$ hdfs dfs -mkdir Project/employment/la
-bash-4.2$ hdfs dfs -mkdir Project/employment/chicago
-bash-4.2$ hdfs dfs -mkdir Project/income
-bash-4.2$ hdfs dfs -mkdir Project/income/la
-bash-4.2$ hdfs dfs -mkdir Project/income/chicago
```

```
-bash-4.2$ hdfs dfs -mkdir Project/population
-bash-4.2$ hdfs dfs -mkdir Project/population/la
-bash-4.2$ hdfs dfs -mkdir Project/population/chicago
```

- Move all files to the HDFS directories respectively:

Chicago education

```
-bash-4.2$ hdfs dfs -put 2020_education_clean_chicago.csv Project/education/chicago
-bash-4.2$ hdfs dfs -put 2021_education_clean_chicago.csv Project/education/chicago
-bash-4.2$ hdfs dfs -put 2022_education_clean_chicago.csv Project/education/chicago
```

- to check files, enter:

```
-bash-4.2$ hdfs dfs -ls Project/education/chicago
```

LA education

```
-bash-4.2$ hdfs dfs -put 2020_education_clean.csv Project/education/la
-bash-4.2$ hdfs dfs -put 2021_education_clean.csv Project/education/la
-bash-4.2$ hdfs dfs -put 2022_education_clean.csv Project/education/la
```

- to check files, enter:

```
-bash-4.2$ hdfs dfs -ls Project/education/la
```

Chicago Employment

```
-bash-4.2$ hdfs dfs -put 2020_employment_clean_chicago.csv Project/employment/chicago  
-bash-4.2$ hdfs dfs -put 2021_employment_clean_chicago.csv Project/employment/chicago  
-bash-4.2$ hdfs dfs -put 2022_employment_clean_chicago.csv Project/employment/chicago
```

- to check files, enter:

```
-bash-4.2$ hdfs dfs -ls Project/employment/chicago
```

LA employment

```
-bash-4.2$ hdfs dfs -put 2020_employment_clean.csv Project/employment/la  
-bash-4.2$ hdfs dfs -put 2021_employment_clean.csv Project/employment/la  
-bash-4.2$ hdfs dfs -put 2022_employment_clean.csv Project/employment/la
```

- to check files, enter:

```
-bash-4.2$ hdfs dfs -ls Project/employment/la
```

Chicago income

```
-bash-4.2$ hdfs dfs -put 2020_income_clean_chicago.csv Project/income/chicago  
-bash-4.2$ hdfs dfs -put 2021_income_clean_chicago.csv Project/income/chicago  
-bash-4.2$ hdfs dfs -put 2022_income_clean_chicago.csv Project/income/chicago
```

- to check files, enter:

```
-bash-4.2$ hdfs dfs -ls Project/income/chicago
```

LA income

```
-bash-4.2$ hdfs dfs -put 2020_income_clean_la.csv Project/income/la  
-bash-4.2$ hdfs dfs -put 2021_income_clean_la.csv Project/income/la  
-bash-4.2$ hdfs dfs -put 2022_income_clean_la.csv Project/income/la
```

- to check files, enter:

```
-bash-4.2$ hdfs dfs -ls Project/income/la
```

Chicago population

```
-bash-4.2$ hdfs dfs -put 2020_population_clean_chicago.csv Project/population/chicago  
-bash-4.2$ hdfs dfs -put 2021_population_clean_chicago.csv Project/population/chicago  
-bash-4.2$ hdfs dfs -put 2022_population_clean_chicago.csv Project/population/chicago
```

- to check files, enter:

```
-bash-4.2$ hdfs dfs -ls Project/population/chicago
```

LA population

```
-bash-4.2$ hdfs dfs -put 2020_population_clean_la.csv Project/population/la  
-bash-4.2$ hdfs dfs -put 2021_population_clean_la.csv Project/population/la  
-bash-4.2$ hdfs dfs -put 2022_population_clean_la.csv Project/population/la
```

- to check files, enter:

```
-bash-4.2$ hdfs dfs -ls Project/population/la
```

Step 6: Creating Hive tables to Query Education, Employment, Income and Population Data

- Enter Hadoop:

```
-bash-4.2$ beeline
```

- Use database:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> use 5200group2
```

Create LA education table for 2020, 2021, 2022 years

```
DROP TABLE IF EXISTS original_la_education;  
  
CREATE EXTERNAL TABLE IF NOT EXISTS original_la_education(  
    geo_id STRING,  
    area_name STRING,  
    population_age_18_24 STRING,  
    high_school_grad_or_higher_18_24 STRING,  
    population_25_34_years STRING,  
    high_school_grad_or_higher_25_34 STRING,  
    population_35_44_years STRING,  
    high_school_grad_or_higher_35_44 STRING,  
    population_45_64_years STRING,  
    high_school_grad_or_higher_45_64 STRING,  
    population_age_65_and_over STRING,  
    high_school_grad_or_higher_65_and_over STRING,  
    `year` STRING  
)  
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'  
WITH SERDEPROPERTIES (  
    "separatorChar" = ",",  
    "quoteChar" = "\"",  
    "escapeChar" = "\\\""  
)  
LOCATION '/user/dkim171/Project/education/la/'
```

```
TBLPROPERTIES ('skip.header.line.count'=2');
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
+-----+
| tab_name |
+-----+
| original_la_education |
+-----+

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select* from
original_la_education LIMIT 10;
```

Create Chicago education table for 2020, 2021, 2022 years

```
DROP TABLE IF EXISTS original_chicago_education;

CREATE EXTERNAL TABLE IF NOT EXISTS original_chicago_education(
    geo_id STRING,
    area_name STRING,
    population_age_18_24 STRING,
    high_school_grad_or_higher_18_24 STRING,
    population_25_34_years STRING,
    high_school_grad_or_higher_25_34 STRING,
    population_35_44_years STRING,
    high_school_grad_or_higher_35_44 STRING,
    population_45_64_years STRING,
    high_school_grad_or_higher_45_64 STRING,
    population_age_65_and_over STRING,
    high_school_grad_or_higher_65_and_over STRING,
    `year` STRING
)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
WITH SERDEPROPERTIES (
    "separatorChar" = ",",
    "quoteChar" = "\"",
    "escapeChar" = "\\"
)
LOCATION '/user/dkim171/Project/education/chicago/'
TBLPROPERTIES ('skip.header.line.count'=2');
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select* from original_chicago_education LIMIT 10;
```

Create LA employment table for 2020, 2021, 2022 years

```
DROP TABLE IF EXISTS original_la_employment;

CREATE EXTERNAL TABLE IF NOT EXISTS original_la_employment(
    geo_id STRING,
    area_name STRING,
    population_age_16_and_over STRING,
    population_age_16_and_over_in_labor_force STRING,
    population_age_16_and_over_in_civ_labor_force_employed STRING,
    population_age_16_and_over_in_civ_labor_force_unemployed STRING,
    population_age_16_and_over_not_in_labor_force STRING,
    own_children_of_the_householder_under_6_years STRING,
    own_children_of_the_householder_under_6_to_17_years STRING,
    `year` STRING
)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
WITH SERDEPROPERTIES (
    "separatorChar" = ",",
    "quoteChar" = "\"",
    "escapeChar" = "\\"
)
LOCATION '/user/dkim171/Project/employment/la/'
TBLPROPERTIES ('skip.header.line.count'=2');
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select* from original_la_employment LIMIT 10;
```

Create Chicago employment table for 2020, 2021, 2022 years

```
DROP TABLE IF EXISTS original_chicago_employment;

CREATE EXTERNAL TABLE IF NOT EXISTS original_chicago_employment(
    geo_id STRING,
    area_name STRING,
    population_age_16_and_over STRING,
    population_age_16_and_over_in_labor_force STRING,
    population_age_16_and_over_in_civ_labor_force_employed STRING,
```

```

population_age_16_and_over_in_civ_labor_force_unemployed STRING,
population_age_16_and_over_not_in_labor_force STRING,
own_children_of_the_householder_under_6_years STRING,
own_children_of_the_householder_under_6_to_17_years STRING,
`year` STRING
)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
WITH SERDEPROPERTIES (
  "separatorChar" = ",",
  "quoteChar" = "\"",
  "escapeChar" = "\\"
)
LOCATION '/user/dkim171/Project/employment/chicago/'
TBLPROPERTIES ('skip.header.line.count'=2');

```

- Check if the table was successfully created:

```

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select* from
original_chicago_employment LIMIT 10;

```

Create LA income table for 2020, 2021, 2022 years

```

DROP TABLE IF EXISTS original_la_income;

CREATE EXTERNAL TABLE IF NOT EXISTS original_la_income(
  geo_id STRING,
  area_name STRING,

  less_than_10k STRING,
  from_10k_to_14999 STRING,
  from_15k_to_24999 STRING,
  from_25k_to_34999 STRING,
  from_35k_to_49999 STRING,
  from_50k_to_74999 STRING,
  from_75k_to_99999 STRING,
  from_100k_to_149999 STRING,
  from_150k_to_199999 STRING,
  from_200k_or_more STRING,
  median_income STRING,

  `year` STRING
)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
WITH SERDEPROPERTIES (
  "separatorChar" = ",",
  "quoteChar" = "\"",

```

```

"escapeChar" = "\\"
)
LOCATION '/user/dkim171/Project/income/la/'
TBLPROPERTIES ('skip.header.line.count'=2');

```

- Check if the table was successfully created:

```

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select* from original_la_income
LIMIT 10;

```

Create Chicago income table for 2020, 2021, 2022 years

```

DROP TABLE IF EXISTS original_chicago_income;

CREATE EXTERNAL TABLE IF NOT EXISTS original_chicago_income(
    geo_id STRING,
    area_name STRING,

    less_than_10k STRING,
    from_10k_to_14999 STRING,
    from_15k_to_24999 STRING,
    from_25k_to_34999 STRING,
    from_35k_to_49999 STRING,
    from_50k_to_74999 STRING,
    from_75k_to_99999 STRING,
    from_100k_to_149999 STRING,
    from_150k_to_199999 STRING,
    from_200k_or_more STRING,
    median_income STRING,

    `year` STRING
)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
WITH SERDEPROPERTIES (
    "separatorChar" = ",",
    "quoteChar" = "\"",
    "escapeChar" = "\\"
)
LOCATION '/user/dkim171/Project/income/chicago/'
TBLPROPERTIES ('skip.header.line.count'=2');

```

- Check if the table was successfully created:

```

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select* from
original_chicago_income LIMIT 10;

```

Create LA Population table for 2020, 2022 years

```
DROP TABLE IF EXISTS original_la_population;

CREATE EXTERNAL TABLE IF NOT EXISTS original_la_population(
    geo_id STRING,
    area_name STRING,
    one_race STRING,
    black_or_African_american STRING,
    merican_indian_and_alaska_native STRING,
    herokee_tribal_grouping STRING,
    chippewa_tribal_grouping STRING,
    navajo_tribal_grouping STRING,
    sioux_tribal_grouping STRING,
    asian STRING,
    indian STRING,
    chinese STRING,
    ilipino STRING,
    japanese STRING,
    korean STRING,
    vietnamese STRING,
    ther_asian STRING,
    native_hawaiian_and_other_pacific_islander STRING,
    chamorro STRING,
    native_hawaiian STRING,
    samoan STRING,
    ther_native_hawaiian_and_other_pacific_pslander STRING,
    race_alone_or_in_combination_with_one_or_more_other_races_total_population STRING,
    hispanic_or_latino_total_population STRING,
    ispanic_or_latino_total_population_of_any_race STRING,
    mexican STRING,
    puerto_rican STRING,
    cuban STRING,
    other_hispanic_or_latin STRING,
    `year` STRING
)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
WITH SERDEPROPERTIES (
    "separatorChar" = ",",
    "quoteChar" = "\"",
    "escapeChar" = "\\"
)
```

```
LOCATION '/user/dkim171/Project/population/la/'  
TBLPROPERTIES ('skip.header.line.count'='2');
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select* from  
original_la_population LIMIT 10;
```

Create Chicago Population table for 2020, 2022 years

```
DROP TABLE IF EXISTS original_chicago_population;  
  
CREATE EXTERNAL TABLE IF NOT EXISTS original_chicago_population(  
    geo_id STRING,  
    area_name STRING,  
  
    one_race STRING,  
    black_or_African_american STRING,  
    merican_indian_and_alaska_native STRING,  
    herokee_tribal_grouping STRING,  
    chippewa_tribal_grouping STRING,  
    navajo_tribal_grouping STRING,  
    sioux_tribal_grouping STRING,  
    asian STRING,  
    indian STRING,  
    chinese STRING,  
    ilipino STRING,  
    japanese STRING,  
    korean STRING,  
    vietnamese STRING,  
    ther_asian STRING,  
    native_hawaiian_and_other_pacific_islander STRING,  
    chamorro STRING,  
    native_hawaiian STRING,  
    samoan STRING,  
    ther_native_hawaiian_and_other_pacific_islander STRING,  
    race_alone_or_in_combination_with_one_or_more_other_races_total_population STRING,  
    hispanic_or_latino_total_population STRING,  
    ispanic_or_latino_total_population_of_any_race STRING,  
    mexican STRING,  
    puerto_rican STRING,  
    cuban STRING,  
    other_hispanic_or_latin STRING,
```

```

`year` STRING
)
ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.OpenCSVSerde'
WITH SERDEPROPERTIES (
  "separatorChar" = ",",
  "quoteChar" = "\"",
  "escapeChar" = "\\"
)
LOCATION '/user/dkim171/Project/population/chicago/'
TBLPROPERTIES ('skip.header.line.count'=2');

```

- Check if the table was successfully created:

```

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select* from
original_chicago_population LIMIT 10;

```

Step 7: Join Hive tables Education, Employment, Income and Population with Crime Data and download for the future analysis

Create Los Angeles crime description table

```

DROP TABLE IF EXISTS crime_description_la;

CREATE TABLE crime_description_la AS
SELECT DISTINCT
  CAST(crm_cd AS INT) AS crm_cd,
  crm_cd_desc
FROM original_la_crime;

```

Join Los Angeles education and crime tables and download final file

```

DROP TABLE IF EXISTS la_education_and_crime;
CREATE TABLE la_education_and_crime AS
SELECT
  SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) AS tract_id,
  c.year,
  CAST(p. primary_type AS INT) AS crm_cd,
  p.count,
  c.population_age_18_24,
  c.high_school_grad_or_higher_18_24,

```

```

c.population_25_34_years,
c.high_school_grad_or_higher_25_34,
c.population_35_44_years,
c.high_school_grad_or_higher_35_44,
c.population_45_64_years,
c.high_school_grad_or_higher_45_64,
c.population_age_65_and_over,
c.high_school_grad_or_higher_65_and_over,
'Los Angeles County' AS area_name
FROM
original_la_education c
LEFT JOIN
la_tractid p
ON
c.year = p.year
AND SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) = p.tract_id
ORDER BY
c.year;

```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from
la_education_and_crime LIMIT 10;
```

Add new column with crime description:

```

DROP TABLE IF EXISTS la_education_with_crime_description;

CREATE TABLE la_education_with_crime_description AS
SELECT t1.*,
REGEXP_REPLACE(t2.crm_cd_desc, ',', '') AS description
FROM
la_education_and_crime t1
LEFT JOIN
crime_description_la t2
ON
t1.crm_cd = t2.crm_cd
ORDER BY t1.year;

```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from
la_education_with_crime_description LIMIT 10;
```

Since the education table has some tractid values and the crime dataset does not, we need to **clean up our output table.**

```
DROP TABLE IF EXISTS final_la_education_and_crime;
```

```
CREATE TABLE final_la_education_and_crime
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE
LOCATION '/user/dkim171/Project/education/final/la'
AS
SELECT *
FROM la_education_with_crime_description
WHERE
crm_cd IS NOT NULL
AND count IS NOT NULL;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from
final_la_education_and_crime LIMIT 10;
```

- Open Linux terminal and check file:

```
hdfs dfs -ls /user/dkim171/Project/education/final/la
```

- Get file from the HDFS to the Linux server:

```
hdfs dfs -get /user/dkim171/Project/education/final/la/000000_0
```

- Exit from the Linux server and enter command in terminal to download file:

```
scp dkim171@129.146.230.230:/home/dkim171/000000_0 ~/Desktop/
```

Join Chicago education and crime tables and download final file

```
DROP TABLE IF EXISTS chicago_education_and_crime;
```

```
CREATE TABLE chicago_education_and_crime AS
SELECT
SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) AS tract_id,
c.year,
p.primary_type,
p.count,
```

```

c.population_age_18_24,
c.high_school_grad_or_higher_18_24,
c.population_25_34_years,
c.high_school_grad_or_higher_25_34,
c.population_35_44_years,
c.high_school_grad_or_higher_35_44,
c.population_45_64_years STRING,
c.high_school_grad_or_higher_45_64,
c.population_age_65_and_over,
c.high_school_grad_or_higher_65_and_over,
'Cook County' AS area_name
FROM
original_chicago_education c
LEFT JOIN
chicago_tractid p
ON
c.year = p.year
AND SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) = p.tract_id
ORDER BY
c.year;

```

- Check if the table was successfully created:

```

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from
final_chicago_education_and_crime LIMIT 10;

```

Clean up our output table.

```

DROP TABLE IF EXISTS final_chicago_education_and_crime;

CREATE TABLE final_chicago_education_and_crime
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE
LOCATION '/user/dkim171/Project/education/final/chicago'

AS
SELECT *
FROM chicago_education_and_crime
WHERE
primary_type IS NOT NULL
AND count IS NOT NULL;

```

- Check if the table was successfully created:

```

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;

```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from final_chicago_education_and_crime LIMIT 10;
```

- Open Linux terminal and check file:

```
hdfs dfs -ls /user/dkim171/Project/education/final/chicago
```

- Remove file from the linux server

```
-bash-4.2$ rm 000000_0
```

- Get new file from the HDFS to the Linux server:

```
hdfs dfs -get /user/dkim171/Project/education/final/000000_0
```

- Exit from the Linux server and enter command in terminal to download file:

```
scp dkim171@129.146.230.230:/home/dkim171/000000_0 ~/Desktop/
```

Join LA employment and crime tables and download final file

```
DROP TABLE IF EXISTS la_employment_and_crime;

CREATE TABLE la_employment_and_crime AS
SELECT
    SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) AS tract_id,
    c.year,
    CAST(p.primary_type AS INT) AS crm_cd,
    p.count,
    c.population_age_16_and_over,
    c.population_age_16_and_over_in_labor_force,
    c.population_age_16_and_over_in_civ_labor_force_employed,
    c.population_age_16_and_over_in_civ_labor_force_unemployed,
    c.population_age_16_and_over_not_in_labor_force,
    c.own_children_of_the_householder_under_6_years,
    c.own_children_of_the_householder_under_6_to_17_years,
    'Los Angeles County' AS area_name
FROM
    original_la_employment c
LEFT JOIN
    la_tractid p
ON
    c.year = p.year
```

```
    AND SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) = p.tract_id  
ORDER BY  
    c.year;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from  
la_employment_and_crime LIMIT 10;
```

Add new column with crime description:

```
DROP TABLE IF EXISTS la_employment_with_crime_description;  
  
CREATE TABLE la_employment_with_crime_description AS  
SELECT t1.*,  
REGEXP_REPLACE(t2.crm_cd_desc, ',', '') AS description  
FROM  
    la_employment_and_crime t1  
LEFT JOIN  
    crime_description_la t2  
ON  
    t1.crm_cd = t2.crm_cd  
ORDER BY t1.year;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from  
la_employment_with_crime_description LIMIT 10;
```

- Clean up our output table.

```
DROP TABLE IF EXISTS final_la_employment_and_crime;  
  
CREATE TABLE final_la_employment_and_crime  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY ','  
STORED AS TEXTFILE  
LOCATION '/user/dkim171/Project/employment/final/la'  
  
AS  
SELECT *  
FROM la_employment_with_crime_description  
WHERE  
   .crm_cd IS NOT NULL  
    AND count IS NOT NULL;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;  
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from final_la_employment_and_crime LIMIT 10;
```

- Open Linux terminal and check file:

```
hdfs dfs -ls /user/dkim171/Project/employment/final/la
```

- Remove file from the Linux server

```
-bash-4.2$ rm 000000_0
```

- Get new file from the HDFS to the Linux server:

```
hdfs dfs -get /user/dkim171/Project/employment/final/la/000000_0
```

- Exit from the Linux server and enter command in terminal to download file:

```
scp dkim171@129.146.230.230:/home/dkim171/000000_0 ~/Desktop/
```

Join Chicago employment and crime tables and download final file

```
DROP TABLE IF EXISTS chicago_employment_and_crime;  
  
CREATE TABLE chicago_employment_and_crime AS  
SELECT  
    SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) AS tract_id,  
    c.year,  
    p.primary_type,  
    p.count,  
  
    c.population_age_16_and_over,  
    c.population_age_16_and_over_in_labor_force,  
    c.population_age_16_and_over_in_civ_labor_force_employed,  
    c.population_age_16_and_over_in_civ_labor_force_unemployed,  
    c.population_age_16_and_over_not_in_labor_force,  
    c.own_children_of_the_householder_under_6_years,  
    c.own_children_of_the_householder_under_6_to_17_years,  
    'Cook County' AS area_name  
FROM  
    original_chicago_employment c  
LEFT JOIN  
    chicago_tractid p  
ON  
    c.year = p.year  
    AND SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) = p.tract_id
```

```
ORDER BY
```

```
c.year;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from  
chicago_employment_and_crime LIMIT 10;
```

- Clean up our output table.

```
DROP TABLE IF EXISTS final_chicago_employment_and_crime;
```

```
CREATE TABLE final_chicago_employment_and_crime  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY ','  
STORED AS TEXTFILE  
LOCATION '/user/dkim171/Project/employment/final/chicago'  
  
AS  
SELECT *  
FROM chicago_employment_and_crime  
WHERE  
primary_type IS NOT NULL  
AND count IS NOT NULL;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from  
final_chicago_employment_and_crime LIMIT 10;
```

- Open Linux terminal and check file:

```
hdfs dfs -ls /user/dkim171/Project/employment/final/chicago
```

- Remove file from the Linux server

```
-bash-4.2$ rm 000000_0
```

- Get new file from the HDFS to the Linux server:

```
hdfs dfs -get /user/dkim171/Project/employment/final/chicago/000000_0
```

- Exit from the Linux server and enter command in terminal to download file:

```
scp dkim171@129.146.230.230:/home/dkim171/000000_0 ~/Desktop/
```

Join LA Income and Crime tables and download final file

```
DROP TABLE IF EXISTS la_income_and_crime;

CREATE TABLE la_income_and_crime AS
SELECT
    SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) AS tract_id,
    c.year,
    CAST(p.primary_type AS INT) AS crm_cd,
    p.count,
    c.less_than_10k,
    c.from_10k_to_14999,
    c.from_15k_to_24999,
    c.from_25k_to_34999,
    c.from_35k_to_49999,
    c.from_50k_to_74999,
    c.from_75k_to_99999,
    c.from_100k_to_149999,
    c.from_150k_to_199999,
    c.from_200k_or_more,
    REGEXP_REPLACE(c.median_income, '250,000\\+', '250000') AS median_income,
    'Los Angeles County' AS area_name
FROM
    original_la_income c
LEFT JOIN
    la_tractid p
ON
    c.year = p.year
    AND SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) = p.tract_id
ORDER BY
    c.year;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from
la_income_and_crime LIMIT 10;
```

- Add new column with crime description:

```
DROP TABLE IF EXISTS la_income_with_crime_description;
```

```
CREATE TABLE la_income_with_crime_description AS
```

```

SELECT t1.*,
REGEXP_REPLACE(t2.crm_cd_desc, ',', '') AS description
FROM
    la_income_and_crime t1
LEFT JOIN
    crime_description_la t2
ON
    t1.crm_cd = t2.crm_cd
ORDER BY t1.year;

```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from
la_income_with_crime_description LIMIT 10;
```

- Clean up our output table.

```
DROP TABLE IF EXISTS final_la_income_and_crime;
```

```

CREATE TABLE final_la_income_and_crime
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE
LOCATION '/user/dkim171/Project/income/final/la'

AS
SELECT *
FROM la_income_with_crime_description
WHERE
    crm_cd IS NOT NULL
    AND count IS NOT NULL;

```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from
final_la_income_and_crime LIMIT 10;
```

- Open Linux terminal and check file:

```
hdfs dfs -ls /user/dkim171/Project/income/final/la
```

- Remove previous file from the Linux server

```
-bash-4.2$ rm 000000_0
```

- Get new file from the HDFS to the Linux server:

```
hdfs dfs -get /user/dkim171/Project/income/final/la/000000_0
```

- Exit from the Linux server and enter command in terminal to download file:

```
scp dkim171@129.146.230.230:/home/dkim171/000000_0 ~/Desktop/
```

Join Chicago Income and Crime tables and download final file

```
DROP TABLE IF EXISTS chicago_income_and_crime;

CREATE TABLE chicago_income_and_crime AS
SELECT
    SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) AS tract_id,
    c.year,
    p.primary_type,
    p.count,
    c.less_than_10k,
    c.from_10k_to_14999,
    c.from_15k_to_24999,
    c.from_25k_to_34999,
    c.from_35k_to_49999,
    c.from_50k_to_74999,
    c.from_75k_to_99999,
    c.from_100k_to_149999,
    c.from_150k_to_199999,
    c.from_200k_or_more,
    c.median_income,
    'Cook County' AS area_name
FROM
    original_chicago_income c
LEFT JOIN
    chicago_tractid p
ON
    c.year = p.year
    AND SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) = p.tract_id
ORDER BY
    c.year;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from  
chicago_income_and_crime LIMIT 10;
```

- Clean up our output table.

```
DROP TABLE IF EXISTS final_chicago_income_and_crime;
```

```
CREATE TABLE final_chicago_income_and_crime  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY ','  
STORED AS TEXTFILE  
LOCATION '/user/dkim171/Project/income/final/chicago'  
  
AS  
SELECT *  
FROM chicago_income_and_crime  
WHERE  
    primary_type IS NOT NULL  
    AND count IS NOT NULL;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from  
final_chicago_income_and_crime LIMIT 10;
```

- Open Linux terminal and check file:

```
hdfs dfs -ls /user/dkim171/Project/income/final/chicago
```

- Remove previous file from the Linux server

```
-bash-4.2$ rm 000000_0
```

- Get new file from the HDFS to the Linux server:

```
hdfs dfs -get /user/dkim171/Project/income/final/chicago/000000_0
```

- Exit from the Linux server and enter command in terminal to download file:

```
scp dkim171@129.146.230.230:/home/dkim171/000000_0 ~/Desktop/
```

Join LA Population and Crime tables and download final file

```

DROP TABLE IF EXISTS la_population_and_crime;

CREATE TABLE la_population_and_crime AS
SELECT
    SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) AS tract_id,
    c.year,
    CAST(p.primary_type AS INT) AS crm_cd,
    p.count,
    c.one_race,
    c.black_or_African_american,
    c.merican_indian_and_alaska_native,
    c.herokee_tribal_grouping,
    c.chippewa_tribal_grouping,
    c.navajo_tribal_grouping,
    c.sioux_tribal_grouping,
    c.asian,
    c.indian,
    c.chinese,
    c.ilipino,
    c.japanese,
    c.korean,
    c.vietnamese,
    c.ther_asian,
    c.native_hawaiian_and_other_pacific_islander,
    c.chamorro,
    c.native_hawaiian,
    c.samoan,
    c.ther_native_hawaiian_and_other_pacific_pslander,
    c.race_alone_or_in_combination_with_one_or_more_other_races_total_population,
    c.hispanic_or_latino_total_population,
    c.ispanic_or_latino_total_population_of_any_race,
    c.mexican,
    c.puerto_rican,
    c.cuban,
    c.other_hispanic_or_latin,
    'Los Angeles County' AS area_name
FROM
    original_la_population c
LEFT JOIN
    la_tractid p
ON
    c.year = p.year
    AND SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) = p.tract_id

```

```
ORDER BY  
c.year;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from  
la_population_and_crime LIMIT 10;
```

- **Add new column with crime description:**

```
DROP TABLE IF EXISTS la_population_with_crime_description;
```

```
CREATE TABLE la_population_with_crime_description AS  
SELECT t1.*,  
REGEXP_REPLACE(t2.crm_cd_desc, ',', '') AS description  
FROM  
la_population_and_crime t1  
LEFT JOIN  
crime_description_la t2  
ON  
t1.crm_cd = t2.crm_cd  
ORDER BY t1.year;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from  
la_population_with_crime_description LIMIT 10;
```

- Clean up our output table.

```
DROP TABLE IF EXISTS final_la_population_and_crime;
```

```
CREATE TABLE final_la_population_and_crime  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY ','  
STORED AS TEXTFILE  
LOCATION '/user/dkim171/Project/population/final/la'  
  
AS  
SELECT *  
FROM la_population_with_crime_description  
WHERE  
crm_cd IS NOT NULL  
AND count IS NOT NULL;
```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from final_la_population_and_crime LIMIT 10;
```

- Open Linux terminal and check file:

```
hdfs dfs -ls /user/dkim171/Project/population/final/la
```

- Remove previous file from the Linux server

```
-bash-4.2$ rm 000000_0
```

- Get new file from the HDFS to the Linux server:

```
hdfs dfs -get /user/dkim171/Project/population/final/la/000000_0
```

- Exit from the Linux server and enter command in terminal to download file:

```
scp dkim171@129.146.230.230:/home/dkim171/000000_0 ~/Desktop/
```

Join Chicago Population and Crime tables and download final file

```
DROP TABLE IF EXISTS chicago_population_and_crime;
CREATE TABLE chicago_population_and_crime AS
SELECT
    SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) AS tract_id,
    c.year,
    p.primary_type,
    p.count,
    c.one_race,
    c.black_or_African_american,
    c.merican_indian_and_alaska_native,
    c.herokee_tribal_grouping,
    c.chippewa_tribal_grouping,
    c.navajo_tribal_grouping,
    c.sioux_tribal_grouping,
    c.asian,
    c.indian,
    c.chinese,
    c.ilipino,
    c.japanese,
    c.korean,
```

```

c.vietnamese,
c.ther_asian,
c.native_hawaiian_and_other_pacific_islander,
c.chamorro,
c.native_hawaiian,
c.samoan,
c.ther_native_hawaiian_and_other_pacific_pslander,
c.race_alone_or_in_combination_with_one_or_more_other_races_total_population,
c.hispanic_or_latino_total_population,
c.ispanic_or_latino_total_population_of_any_race,
c.mexican,
c.puerto_rican,
c.cuban,
c.other_hispanic_or_latin,
'Cook County' AS area_name
FROM
original_chicago_population c
LEFT JOIN
chicago_tractid p
ON
c.year = p.year
AND SUBSTRING(c.geo_id, LOCATE('US', c.geo_id) + 2) = p.tract_id
ORDER BY
c.year;

```

- Check if the table was successfully created:

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> show tables;
```

```
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> select * from
chicago_population_and_crime LIMIT 10;
```

- Clean up our output table.

```
DROP TABLE IF EXISTS final_chicago_population_and_crime;
```

```

CREATE TABLE final_chicago_population_and_crime
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE
LOCATION '/user/dkim171/Project/population/final/chicago'

AS
SELECT *
FROM chicago_population_and_crime
WHERE
primary_type IS NOT NULL
AND count IS NOT NULL;

```

- Check if the table was successfully created:
0: jdbc:hive2://bigdaiun0.sub03291929060.trai> **show tables;**

0: jdbc:hive2://bigdaiun0.sub03291929060.trai> **select * from final_chicago_population_and_crime LIMIT 10;**

- Open Linux terminal and check file:
hdfs dfs -ls /user/dkim171/Project/population/final/chicago

- Remove previous file from the Linux server
-bash-4.2\$ **rm 000000_0**

- Get new file from the HDFS to the Linux server:
hdfs dfs -get /user/dkim171/Project/population/final/chicago/000000_0

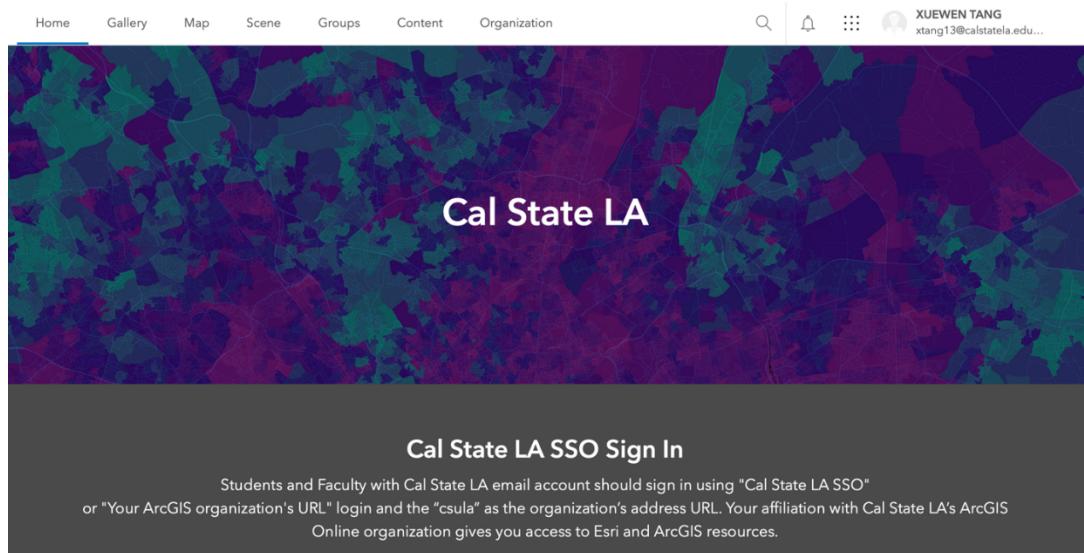
- Exit from the Linux server and enter command in terminal to download file:
scp **dkim171@129.146.230.230:/home/dkim171/000000_0** ~/Desktop/

STEP 8: ArcGIS Online Tutorial

Access Your Account

1. Login with Your Organization Account:
 - Navigate to **ArcGIS Online**: <https://csula.maps.arcgis.com/home/index.html>
 - Select “**Cal State LA SSO**” or “**Your ArcGIS organization’s URL**”.
 - Enter “**csula**” as the organization’s address URL.

Example username: **xtang13@calstatela.edu**. You need to **change the user domain name to your own**.

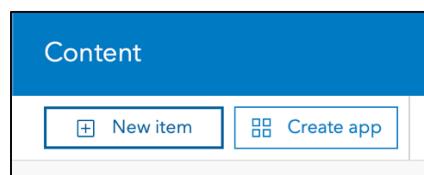


Create feature layer for Chicago and LA:

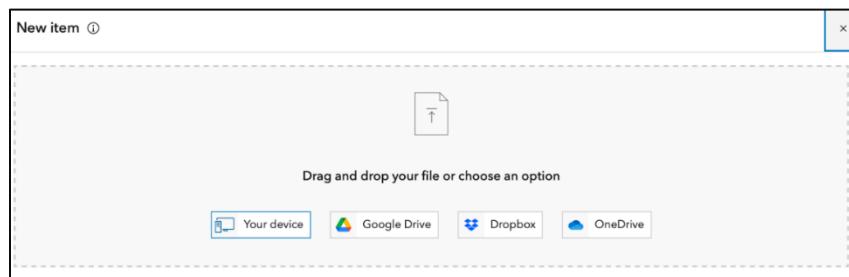
1. Open the “Content” section from the navigation bar.



2. Click “New Item” Button.



3. Upload the aggregate file “crime_chicago_aggregated.zip” in the map data folder. Then click “Next”, Title keeps the same as the file name. Tags: crime, Chicago and click “Save”.



Name	Date Modified	Size	Kind
crime_chicago_aggregated.cpg	Nov 5, 2024 at 8:57 PM	5 bytes	Document
crime_chicago_aggregated.dbf	Nov 5, 2024 at 8:57 PM	469 KB	Document
crime_chicago_aggregated.prj	Nov 5, 2024 at 8:57 PM	145 bytes	Document
crime_chicago_aggregated.shp	Nov 5, 2024 at 8:57 PM	8 MB	Document
crime_chicago_aggregated.shx	Nov 5, 2024 at 8:57 PM	32 KB	Document
crime_chicago_aggregated.zip	Nov 5, 2024 at 8:57 PM	8.5 MB	ZIP archive
crime_la_aggregated.cpg	Nov 9, 2024 at 3:11PM	5 bytes	Document
crime_la_aggregated.dbf	Nov 9, 2024 at 3:11PM	912 KB	Document
crime_la_aggregated.prj	Nov 9, 2024 at 3:11PM	145 bytes	Document
crime_la_aggregated.shp	Nov 9, 2024 at 3:11PM	6.4 MB	Document
crime_la_aggregated.shx	Nov 9, 2024 at 3:11PM	37 KB	Document
crime_la_aggregated.zip	Nov 9, 2024 at 3:11PM	7.4 MB	ZIP archive

New item

File
crime_chicago_aggregated.zip

File type

Shapefile
A vector data storage format for storing the location, shape, and attributes of geographic features. A shapefile is stored in a set of related files and contains one feature class.

How would you like to add this file?

Add crime_chicago_aggregated.zip and create a hosted feature layer
Add the shapefile and publish as a hosted feature layer that can be added to a map.

Add crime_chicago_aggregated.zip only
Add shapefile without publishing. File can be shared and downloaded by others or published at a later date.

Back Cancel Next

New item

File
crime_chicago_aggregated.zip

Title

Folder

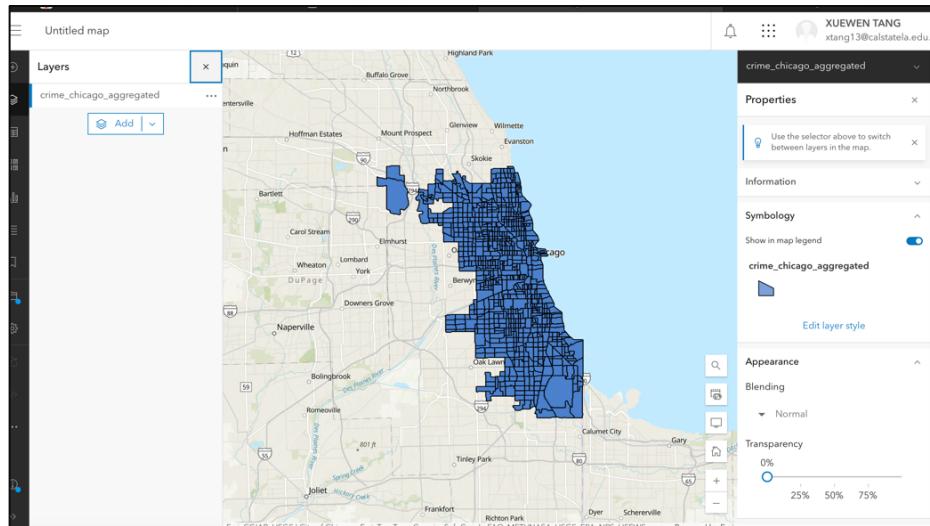
Tags

Summary

0/2048

Back Cancel Save

Then wait for the page processing the file, you will get the map as below:



4. Repeat above steps to create the feature layer for LA.

Then you can get two feature layers ready for Chicago and LA in my content.

<input type="checkbox"/> crime_la_aggregated	Feature Layer (hosted)	Nov 9, 2024		Preview	...
<input type="checkbox"/> crime_la_aggregated	Shapefile	Nov 9, 2024		Preview	...
<input type="checkbox"/> crime_chicago_aggregated	Feature Layer (hosted)	Oct 31, 2024		Preview	...
<input type="checkbox"/> crime_chicago_aggregated	Shapefile	Oct 31, 2024		Preview	...

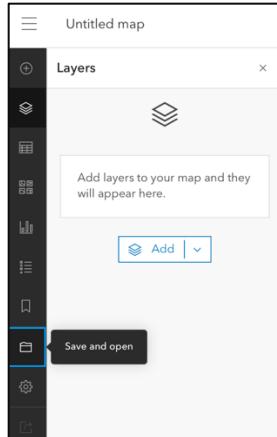
Create a Web Map

1. Open the “Map” section from the navigation bar.



2. Save the map:

- Click “Save As”.



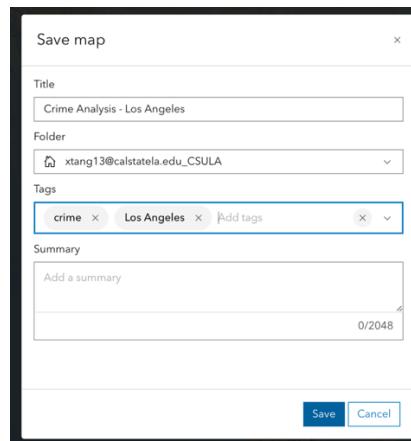
- Complete the details:

Title: Crime Analysis – Los Angeles

Folder: Default (no changes needed)

Tags: Crime, Los Angeles

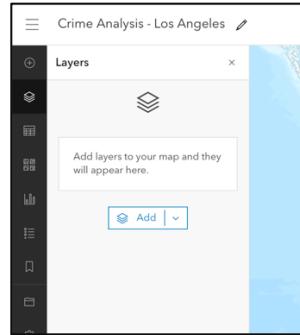
Click “Save”.



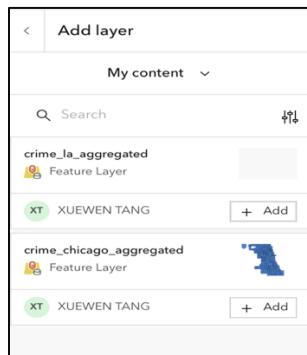
3. Add and Configure Layers

- **Add a Layer:**

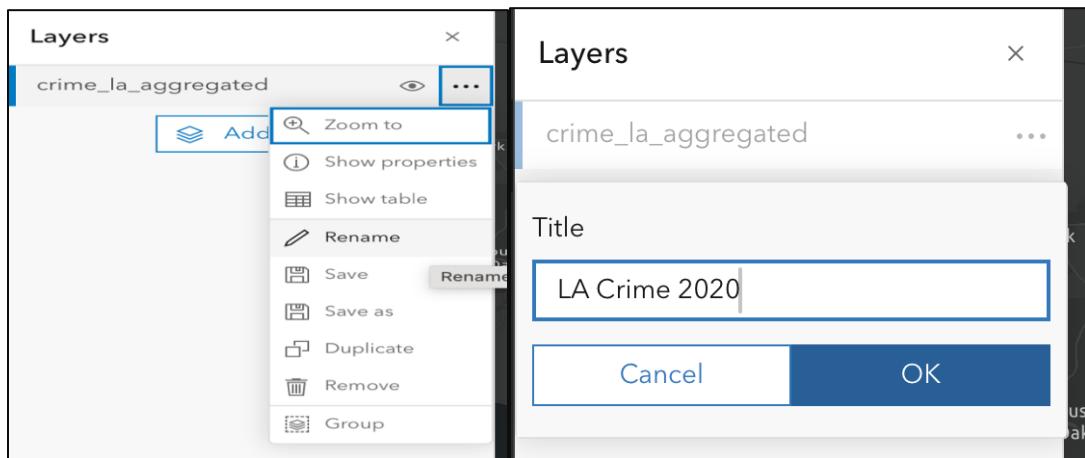
-Click “Add”.



-From **My Content**, select **crime_la_aggregated** and click "+ Add".

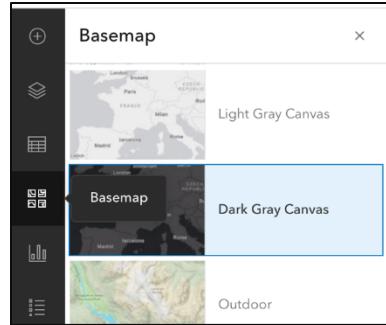


-Rename the layer to “**LA Crime 2020**” for clarity.

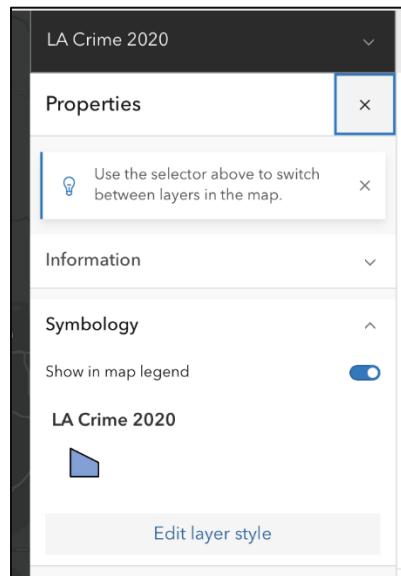


- **Set Basemap:**

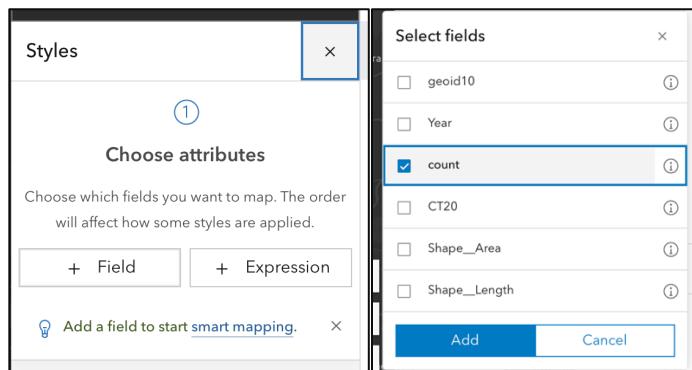
-Click “**Basemap**” and select “**Dark Grey Canvas**”.



- **Edit Layer Style:**
- Click Properties and click “Edit Layer Style”.

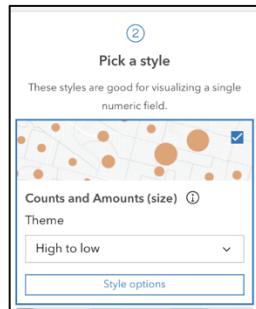


- Choose attributes:
- Click ”+ Field” and select count for the legend.

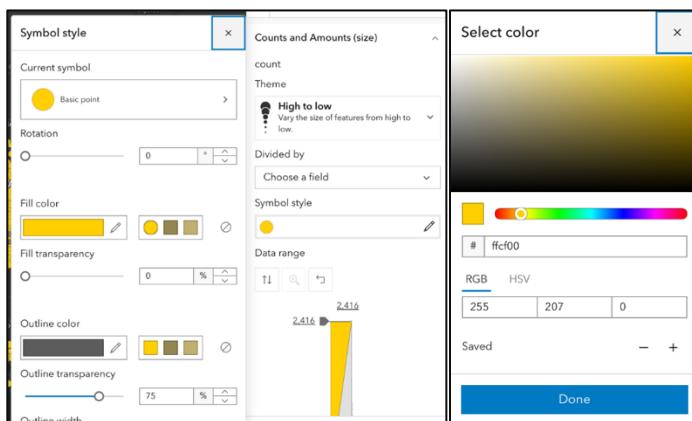
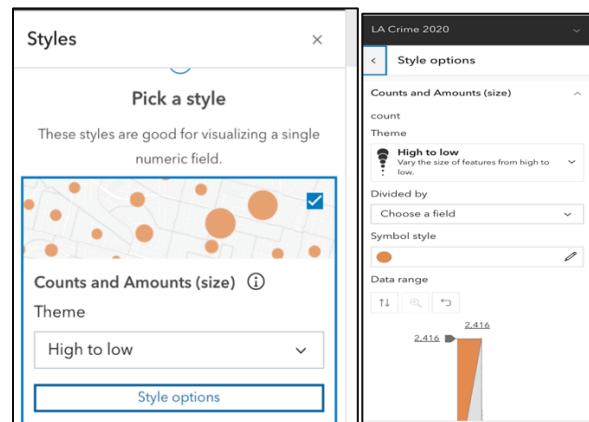


- Pick a style:

Use “Counts and Amounts (size)”.



- Adjust options: Click “Style options”--“Symbol style” edit.



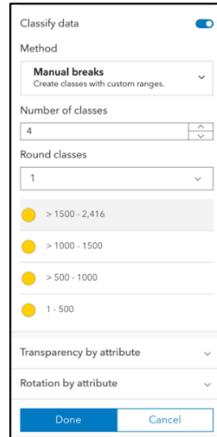
Bubble size: Set from 4 to 40.

Bubble color: **Fill:** #ffcf00. **Outline:** #2b2406.

- Modify data range: Open Classify data

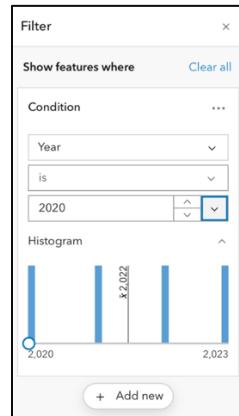
Use **Manual breaks** method with 4 classes and round classes to 1.

Set ranges as: 1–500, 500–1000, 1000–1500, 1500–maximum.



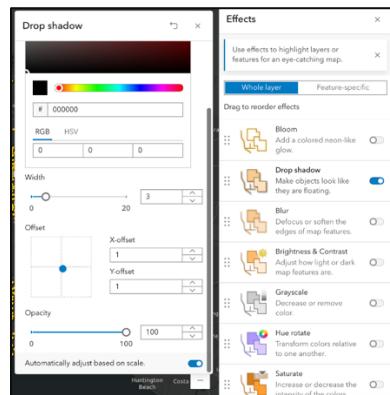
4. Apply Filters:

- Click “Filter” and set Year = 2020.



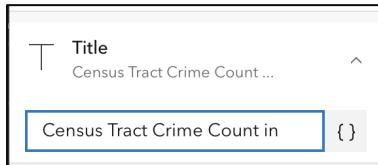
5. Enable Effects:

- Click “Effects” and enable Drop shadow:
- X-offset = 1, Y-offset = 1.

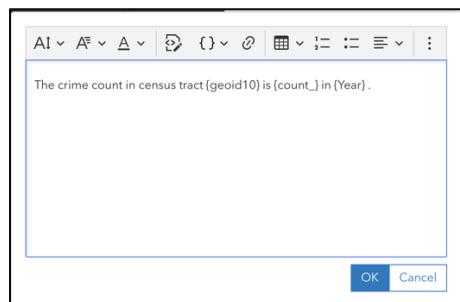


6. Customize Pop-Ups:

- Set the title: Census Tract Crime in {Year}. set “Year” as variable



- Delete default content:
- Add text: The crime count in census tract {geoid10} is {count_} in {Year}.

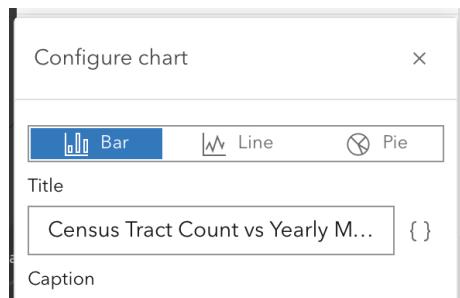


- Add a bar chart media: select “Bar”

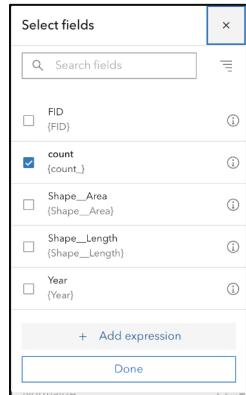
Title: Comparison Analysis.

- Edit media content:

Title: Census Tract Count vs Yearly Maximum



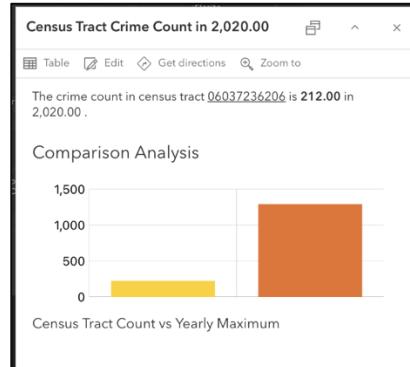
- Select fields: Count



- Add Expression: Name: Year Maximum. Formula: Max(\$layer, "count_").

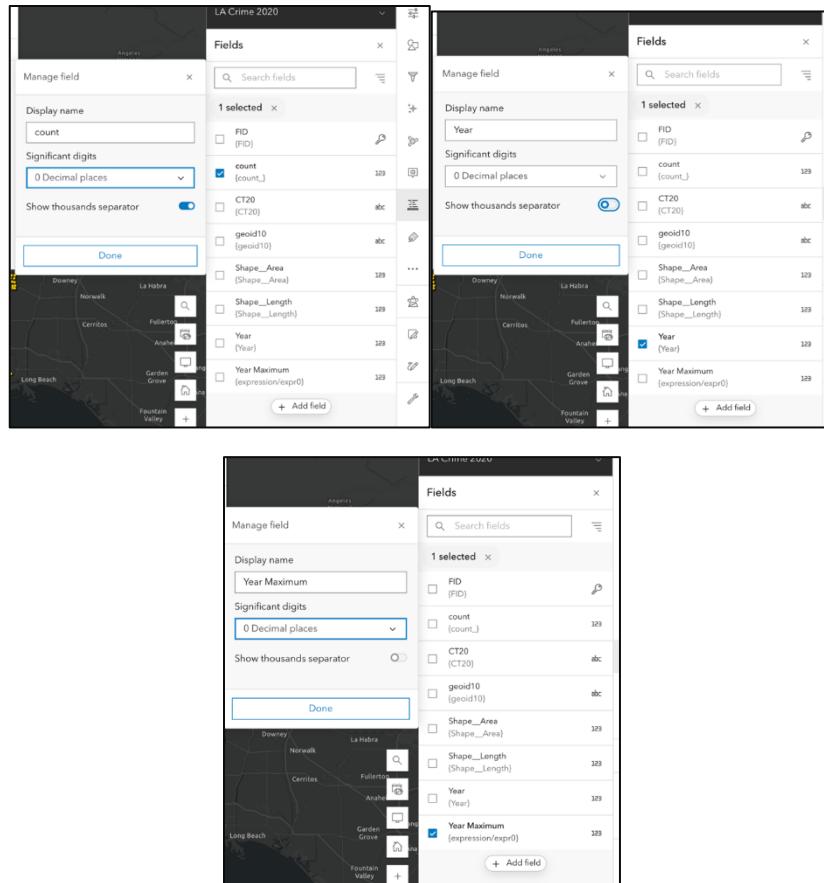


Then you can get the pop-up as below and you can find the display problem in count and year, there is decimal in the number.

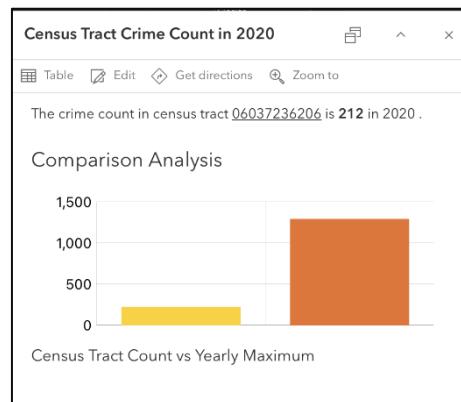


7. Format Fields:

- Click “Fields”, select “count”, change decimal to 0 and show thousands separator open. Select “Year” and “Year Maximum”, change decimal to 0 and close thousands separator.



Then you can see the pop-up changed.

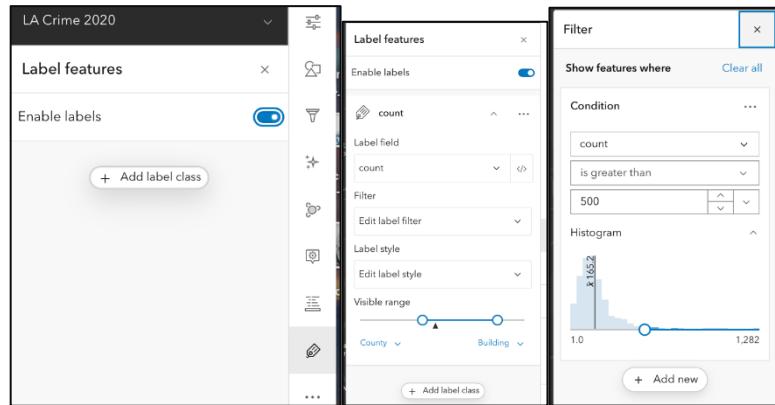


8. Label Layers:

Click “Labels” to add label class. Click “+Add label class”, change label field as “count” and make visible range from County to Building.

- Add labels: **Field:** count. **Visible Range:** From County to Building.

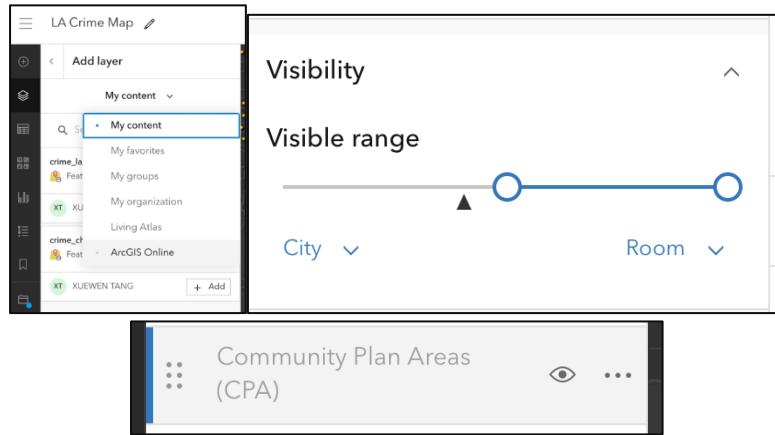
- Filter labels: Edit label filter: click “Add new”, set condition “count” “is greater than””500”



Add Supporting Layers

1. Community Plan Areas:

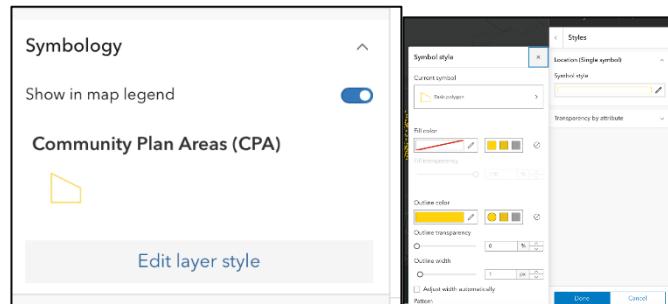
- Search “Community Plan Areas (CPA)” in ArcGIS Online.
- Add the layer and move it to the bottom.
- Open Properties, Set visibility range: City to Room.



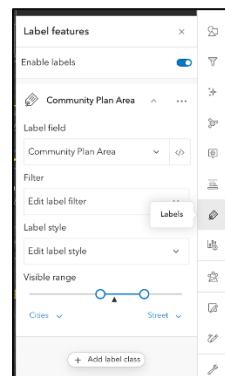
- Style:

Edit Layer style: pick the Location (single symbol) style.

Open style options: edit the symbol style, set the Fill color to “no color” and set the outline color to yellow.



- Label Features: Click “Labels”. Select Label field as “Community Plan Area”, set the visible range from Cities to Street.

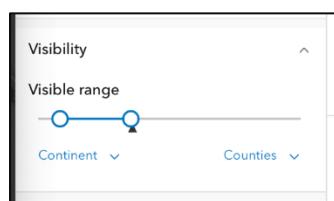


2. LA City Council Districts:

- Search “**LA City Council Districts (Adopted 2021)**” in **ArcGIS Online**.
- Add the layer and move it to the bottom.

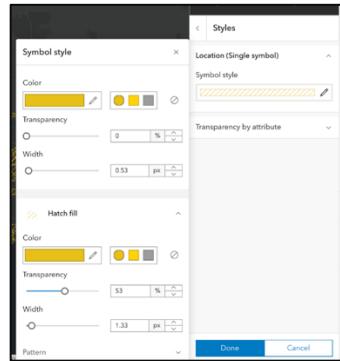


- Open Properties , set visibility range: Continent to Counties.

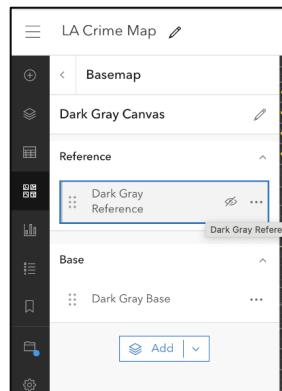


- Style: Edit Layer style: pick the Location (single symbol) style.

Open style options: edit the symbol style, set the **Solid stroke** to “yellow” and select the **Hatch fill** and set the color to “yellow”. **Transparency** set to 53%.

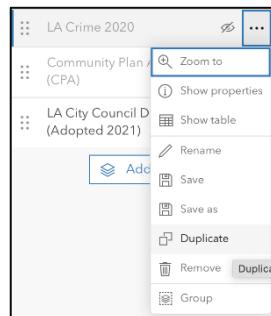


Close Basemap Reference



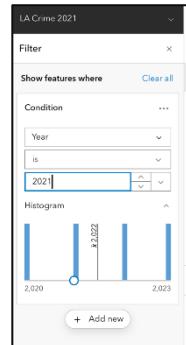
Duplicate for Additional Years

1. Duplicate **LA Crime 2020**:
 - Click “...” next to the layer and select “**Duplicate**”.



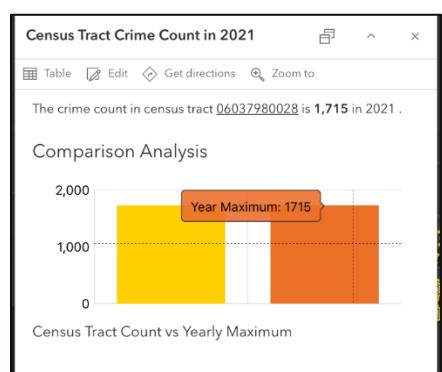
2. Rename and Filter:
 - Rename the layer to **LA Crime [Year]** (e.g., 2021).

- Update the filter to **Year = [Year]**.

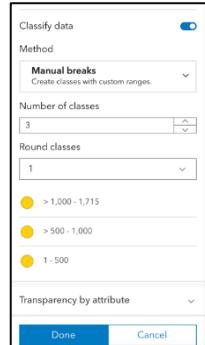


3. Adjust Style Options: Click “Properties”-“Edit Layer style”-“Style Options”.
- Update the maximum data range based on the year’s max value.

- For 2021, as example :Select any bubble to see the pop-up, you can find the maximum count number in this year, for year 2021 the maximum number is 1715.
- Max: 1715.
- Classes: 3.

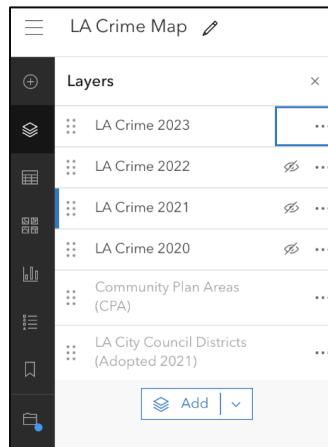


Ranges: 1–500, 500–1000, 1000–1715.



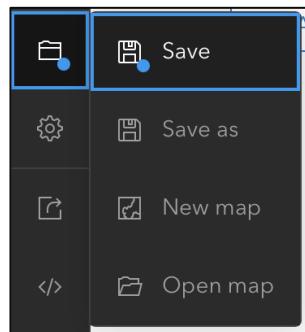
Repeat for 2022 and 2023.

Then you will finish the LA Crime Map as below:



Save and Access Your Map

1. Click “Save” to preserve changes.



2. Find your web map in “My Content”.

Scene Groups Content Organization

My content My favorites My groups My organization Living Atlas

Search All my content

1-7 of 7

Title Modified

LA Crime Map Web Map Nov 17, 2024 Preview ...

Repeat the create web map step of LA above, create the Chicago Crime Map.

Chicago Crime Map Web Map Nov 17, 2024 Preview ...

Create a Map Story

1. Go back to : <https://csula.maps.arcgis.com/home/index.html>
2. Go to “Content” on the navigation bar.

Home Gallery Map Scene Groups Content Organization

3. Click the extend Select “ ArcGIS StoryMaps”

ArcGIS Online ArcGIS StoryMaps ArcGIS for Power BI

Business Analyst Community Analyst Dashboards

Data Pipelines Dev Experience Builder

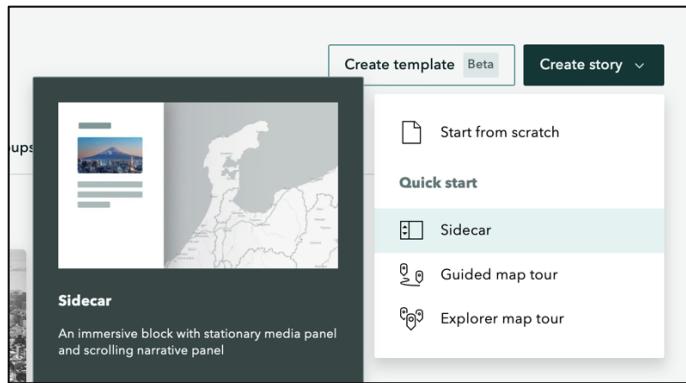
Field Maps Designer Field Maps Data Quality Access

Hub Insights Instant Apps

4. Click “Create Story”

Stories Create template Create story

5. You can select the template or a blank scratch.



6. Edit your Story title, subtitle. You can upload cover image or video.

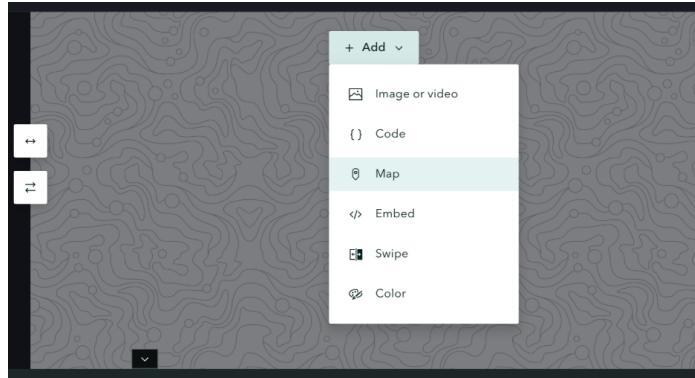
The screenshot shows the Storyline editor with a blank 'Story title' section. A placeholder text 'Start with a short introduction or subtitle (optional)' is visible. Below the title area, it says 'XUEWEN TANG' and 'Draft'. There is a 'Save' button at the top and a 'Publish' button at the top right.

The screenshot shows the published version of the story. The title 'Crime Analysis' is displayed. The subtitle reads: 'Chicago has the highest crime level in the United States during the Pandemic. Has anything changed after the Pandemic?'. Below the text, it says 'CIS 5200 Group 2' and 'November 16, 2024'. To the right of the text is a large, grainy black and white photograph of a city skyline, likely Chicago, viewed through binoculars.

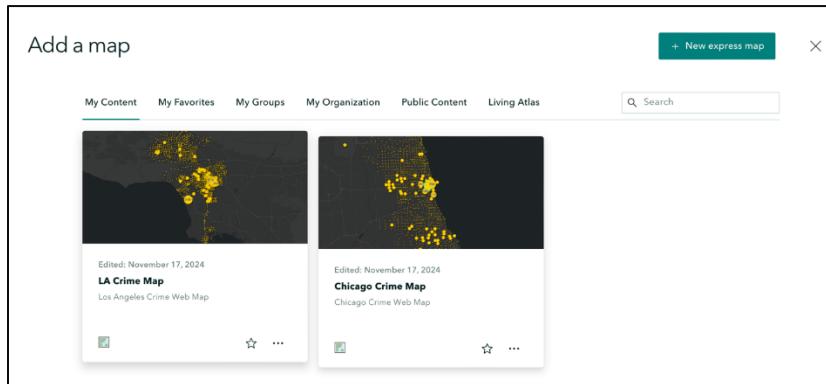
7. Add sidebar
8. Heading: Study Area – Chicago. Double click the text and change paragraph to Heading 1.

The screenshot shows two parts of the Storyline interface. On the left is the sidebar editor with a sidebar menu containing options like Title, Media, Map, Image, etc., and a preview of the sidebar block. On the right is the rich text editor for the 'Study Area - Chicago' section. The heading 'Study Area - Chicago' is selected and highlighted in green. The text area contains nested headings: 'Heading 1' (selected), 'Heading 2', 'Heading 3', and 'Heading 4'.

9. Click "+Add" at the right side, select "Map".



10. You will see the web map we created; it can be selected as below: choose Chicago Crime Map



11. Set the Layer: City of Chicago's Boundaries, as visible, set other layers as invisible.

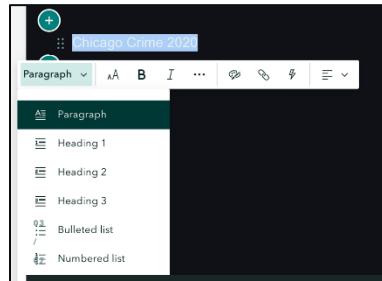


12. Click "New slide" in the bottom Sidecar. Add new slide page.

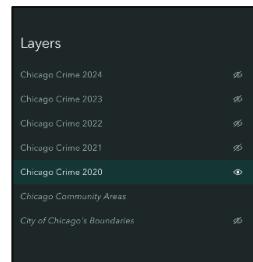


13. Repeat the same step as the first page:

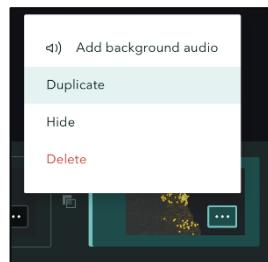
Title: Chicago Crime 2020; Paragraph



14. Add Map: Set the Layer—Chicago Crime 2020; Chicago Community Areas, as visible, set other layers as invisible.

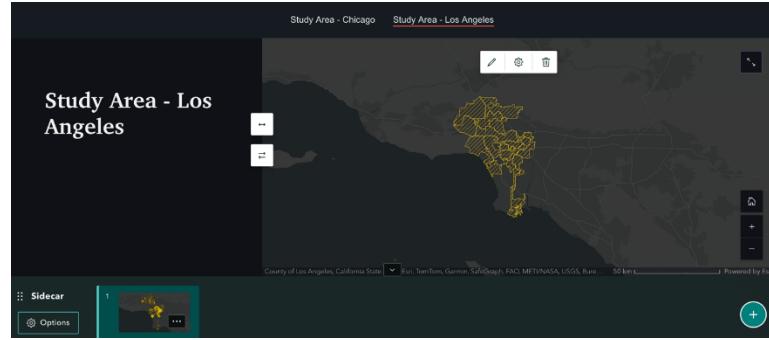


15. You can use Duplicate to add Chicago Crime 2021/2022/2023/2024.
Remember to change the year for each slides.



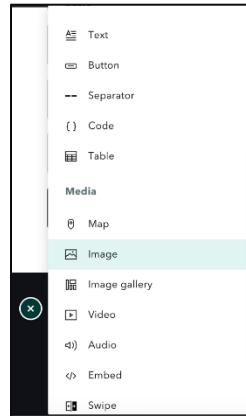
16. Add sidebar for Los Angeles:

Heading: Study Area – Los Angeles. Double click the text and change paragraph to Heading 1.



Add slide and repeat the same steps as Chicago to add slides for year 2020/2021/2022/2023.

17. You can export PPT slides as PNG image, import the PPT slide to the map story.



Step 9: Visual Analysis of Los Angeles and Chicago Crime

Topic: Employment and Crime

Download csv files from Hadoop, in this case, `chicago_employment_and_crime` and `la_education_and_crime` as csv files to local computer.

1. Open PowerBI → Home tab → Get Data → text/CSV → (upload all downloaded tables to PowerBI) shown as below

The screenshot shows the Power BI desktop application. On the left, the 'Home' ribbon is selected. In the center, the 'Data' view is open, displaying a search bar and a list of datasets. The list includes:

- > chicago_education_crime_data
- > chicago_employment_crime_data
- > la_education_crime_data
- > la_employment_crime_data

For Chicago's employment and crime dataset

2. Select stacked area chart as visualization.

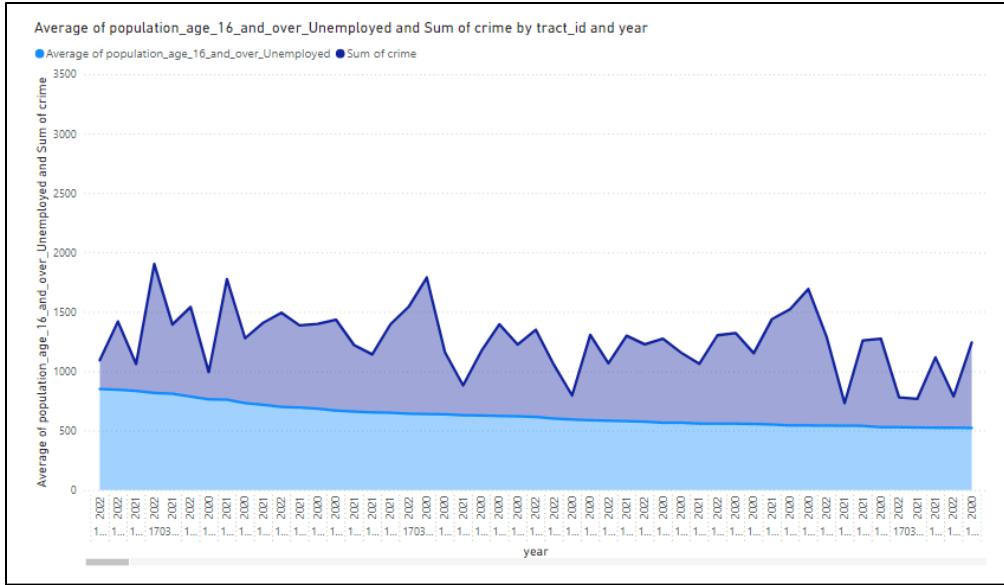
The screenshot shows the 'Visual settings' pane for the 'chicago_employment_crime_data' dataset. The pane lists various fields with checkboxes for selection. The following fields are checked and highlighted with red boxes:

- \sum (U)population_age_16_and_over_in_civ_labor_force_unemployed
- \sum count
- \sum tract_id
- \sum year

Select

- tract_id and year for X-axis,
- [sum] count and [average] (U)population_age_16_and_over_in_civ_labor_force_unemployed for Y-axis

Results should be similar or equivalent as chart below.



For Los Angeles's employment and crime dataset

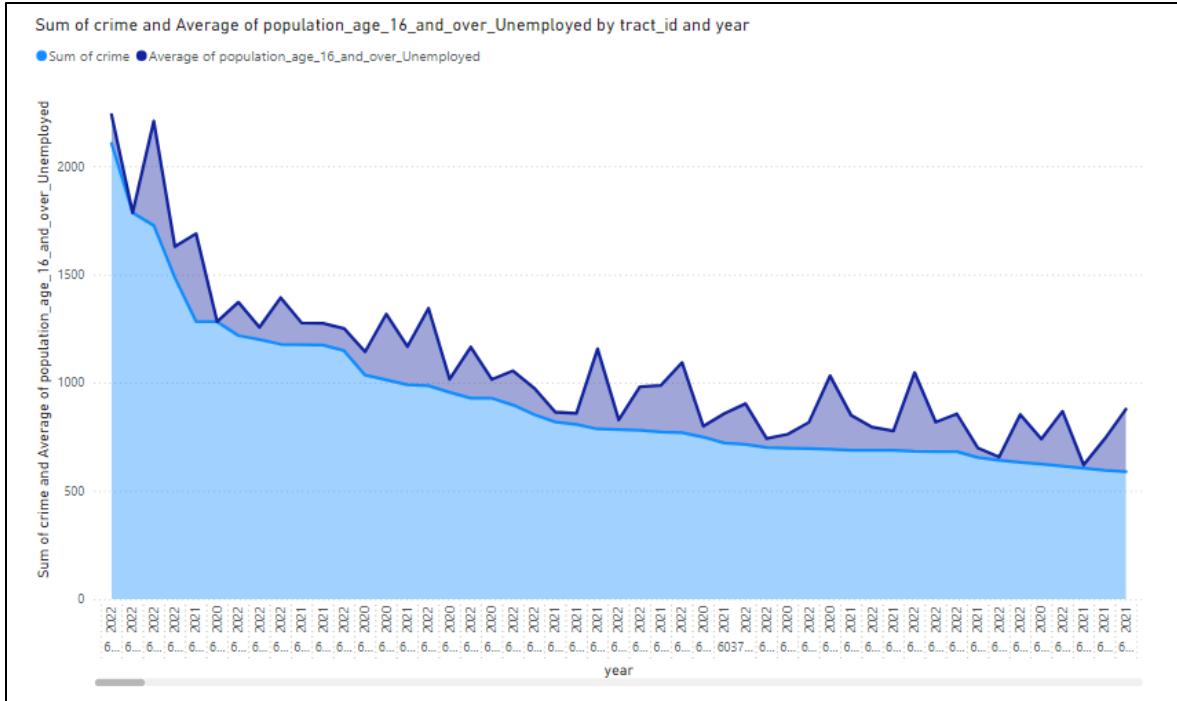
3. Select line chart as visualization

The screenshot shows a dropdown menu for selecting metrics from the 'la_employment_crime_data' dataset. The following metrics are selected and highlighted with red boxes:

- $\sum (U)\text{population_age_16_and_over_in_civ_labor_force_unemployed}$
- $\sum \text{count}$
- $\sum \text{tract_id}$
- $\sum \text{year}$

Select

- tract_id and year for X-axis,
 - [sum] count and [average]
- (U)population_age_16_and_over_in_civ_labor_force_unemployed for Y-axis



Topic: Education and Crime For Chicago's education and crime dataset

1. Select line chart as visualization.

Under employment and crime dataset, select all fields listed in the picture below.

✓ chicago_education_crime_data

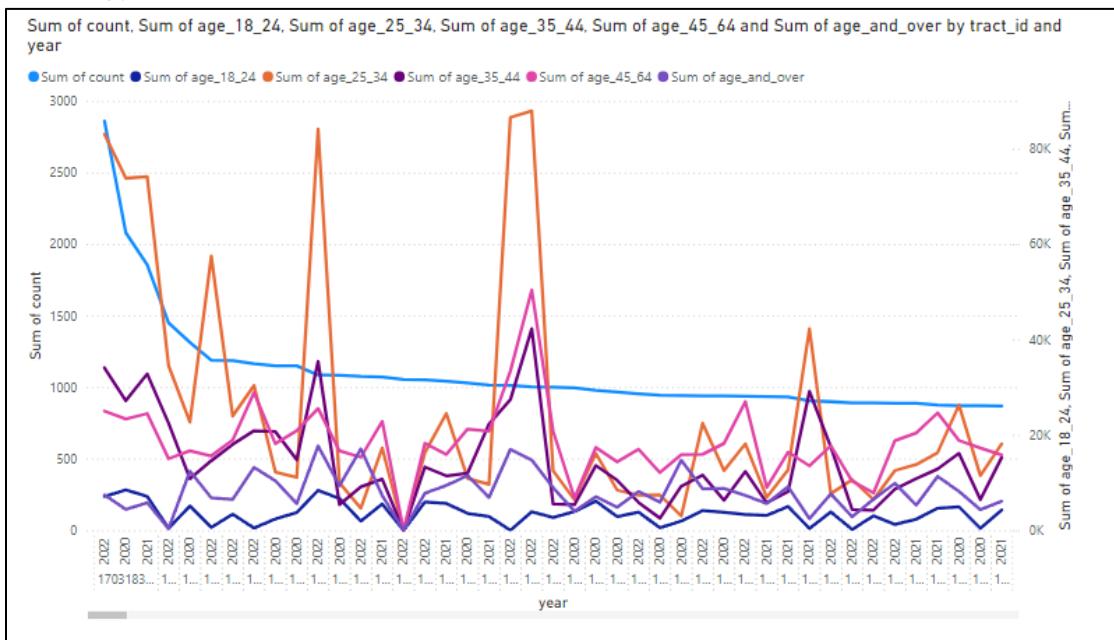
- area_name
- Σ count
- Σ high_school_grad_or_higher_18_24
- Σ high_school_grad_or_higher_25_34
- Σ high_school_grad_or_higher_35_44
- Σ high_school_grad_or_higher_45_64
- Σ high_school_grad_or_higher_65_and_over
- Σ population_25_34_years
- Σ population_35_44_years
- Σ population_age_18_24
- Σ population_age_65_and_over
- primary_type
- Σ string
- Σ tract_id
- Σ year

Select

- tract_id and year for X-axis,

- sum of count for Y-axis,
- [sum] high_school_grad_or_higher_18_24 and the rest of the checkboxes for Secondary Y-axis

As a result, the chart should show up in this format shown below. Apply the same methodology to LA's education and crime dataset.



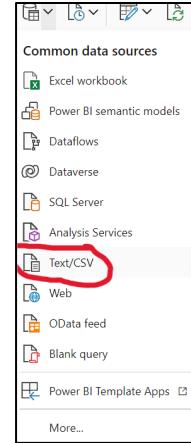
Apply the method for la_education_crime_data

la_education_crime_data

- area_name
- Σ count
- Σ high_school_grad_or_higher_18_24
- Σ high_school_grad_or_higher_25_34
- Σ high_school_grad_or_higher_35_44
- Σ high_school_grad_or_higher_45_64
- Σ high_school_grad_or_higher_65_and_over
- Σ population_25_34_years
- Σ population_35_44_years
- Σ population_age_18_24
- Σ population_age_65_and_over
- Σ primary_type
- Σ string
- Σ tract id

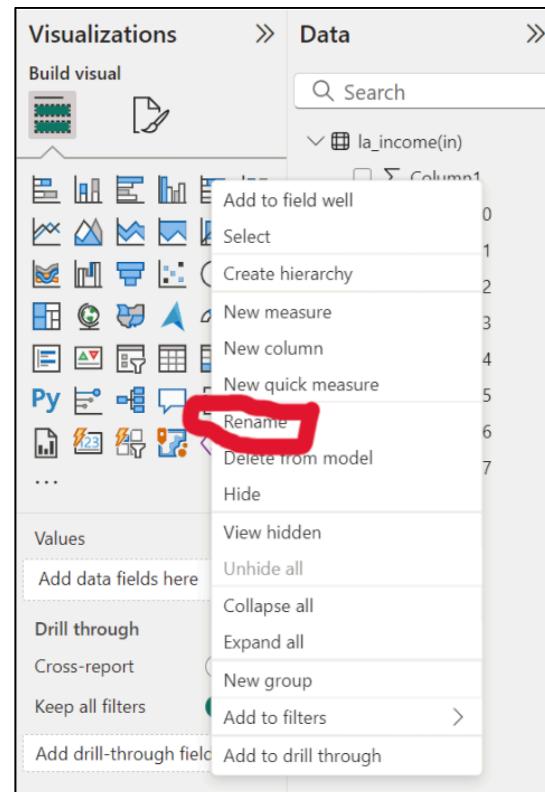
Topic: LA Income and Crime

Step 1: Set Up Power BI Open Power BI Desktop. Click "Get Data" → Select "Text/CSV". Browse to your dataset file and load it into Power BI.



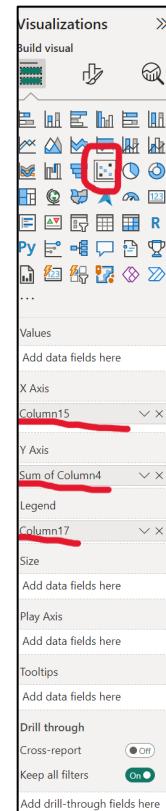
Step 2: Load and Prepare the Data
simply click load
step 3: label all columns
correspondingly:

COLUMN 1 = tract_id
COLUMN 2 = year
COLUMN 3 = crm_cd
COLUMN 4 = count
COLUMN 5 = less_than_10k
COLUMN 6 = from_10k_to_14999
COLUMN 7 = from_15k_to_24999
COLUMN 8 = from_25k_to_34999
COLUMN 9 = from_35k_to_49999
COLUMN 10 = from_50k_to_74999
COLUMN 11 = from_75k_to_99999
COLUMN 12 = from_100k_to_149999
COLUMN 13 = from_150k_to_199999
COLUMN 14 = from_200k_or_more
COLUMN 15 = median_income
COLUMN 16 = area_name
COLUMN 17 = description

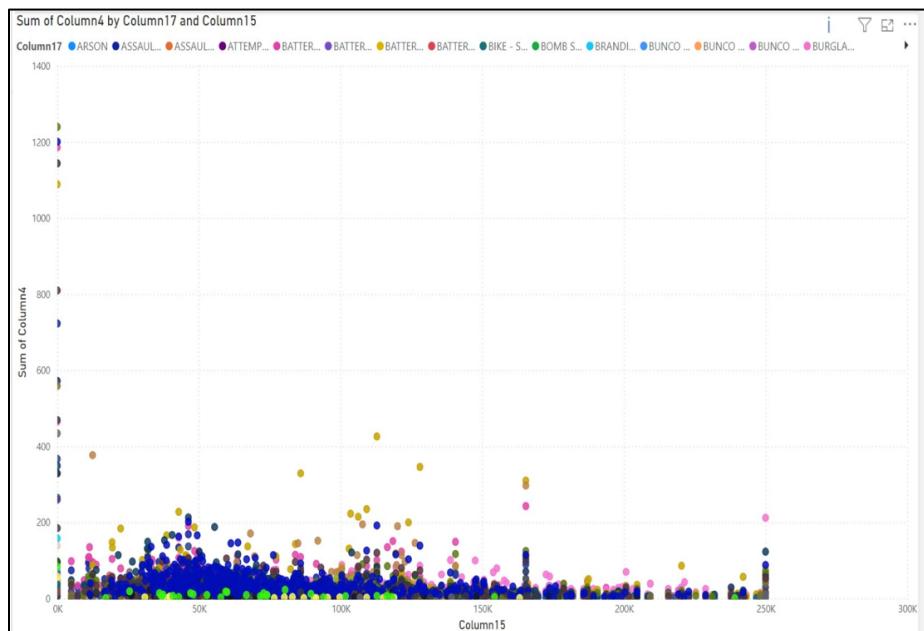


Step 4:

Create the Scatter Chart In the Report View, click on Scatter Chart from the Visualizations pane. Assign the fields to the scatter chart: X-Axis: Drag Column 15 (median_income). Y-Axis: Drag Column 4 (count). Legend: Drag Column 17 (description).

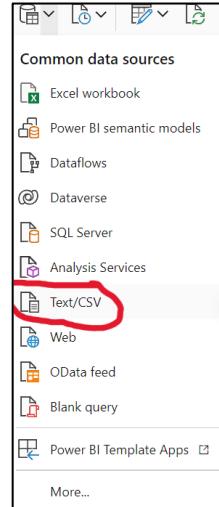


RESULTS:



Treemap LA INCOME:

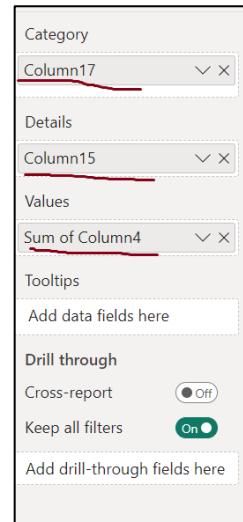
Step 1: Set Up Power BI Open Power BI Desktop. Click "Get Data" → Select "Text/CSV". Locate and import your dataset.



Step 2: Add a Treemap Visualization Navigate to the Report View in Power BI Desktop. From the Visualizations Pane, click on the Treemap icon to add it to your canvas.



Step 3: Configure the Treemap Assign the following fields to the treemap:
Category: Drag Column 17: (description) into the Category field well.
Details: Column 15:
Values (Details): Drag Column 4 (median_income) into the Values field well.
The treemap will now show categories (from description) sized based on their corresponding median income values.



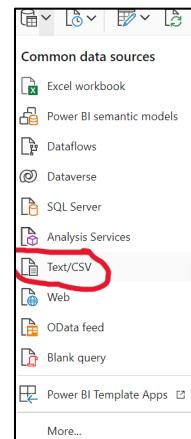
Topic: CHICAGO INCOME and CRIME:

STEP 1: open PowerBI

STEP 2: create a blank document

STEP 3: in the GET DATA module
under home
click text/csv

STEP 4: in the data preview window
click load



STEP 5:

Under Data Visualization Click the
Scatter chart icon

STEP 6: Column 15 in the "Data"
Section represents "median income
data"

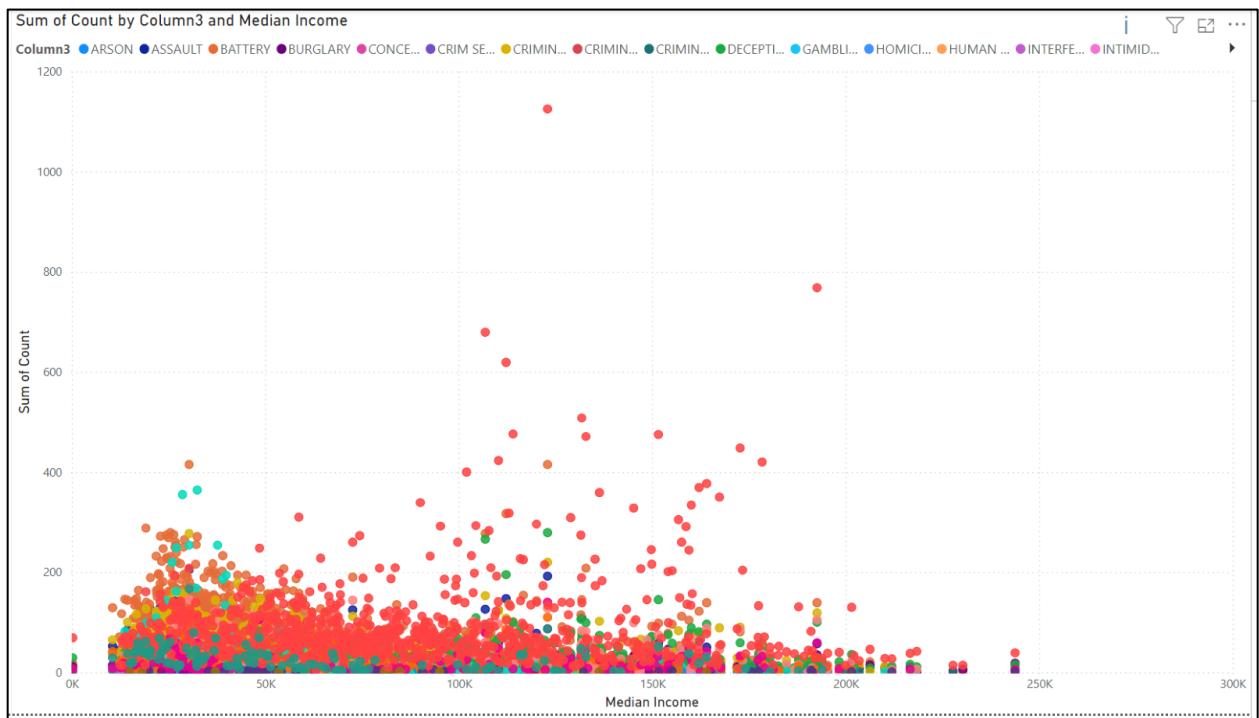
STEP 7: Drag Column 15 to the x-
axis

STEP 8: Drag Column 4 into the Y
axis portion this will represent the
"Count" (frequency of crimes)

The figure consists of three screenshots. The first is the 'Visualizations' pane showing various chart icons, with the scatter plot icon highlighted. The second is the 'X Axis' settings panel where 'Median Income' is listed under 'Median Income'. The third is a context menu for 'Median Income' in the X-axis settings, with the 'Don't summarize' option highlighted.

STEP 9: Under the x column click the drop down menu and put "do not Summarize"
This scatter plot show the positive correlation between income and the frequency
of crime within Chicago. It signifies that everything from 100k and below has a
significantly higher crime rate frequency.

RESULTS:



TREEMAP MAP CHICAGO INCOME/CRIME:

STEP 1: Open a new PowerBI Blank Document

STEP 2: in the GET DATA

module under home

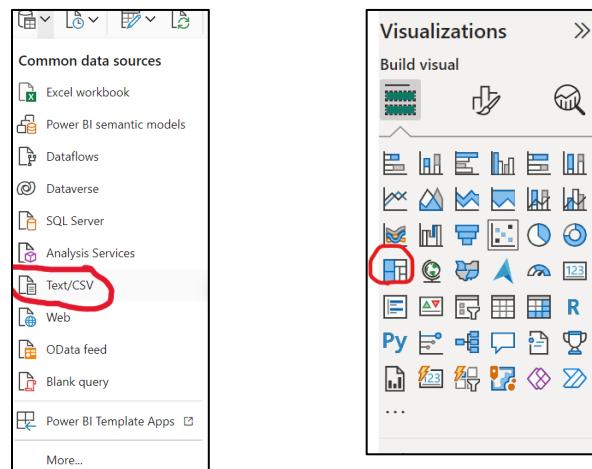
click text/csv

STEP 3: In the data preview

window click load

STEP 4: In the Visualizations

menu click "TreeMap" .

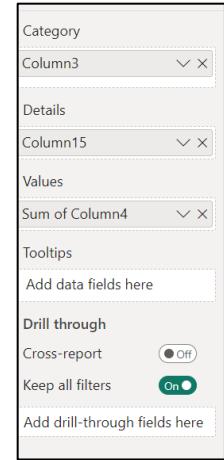


STEP 5:

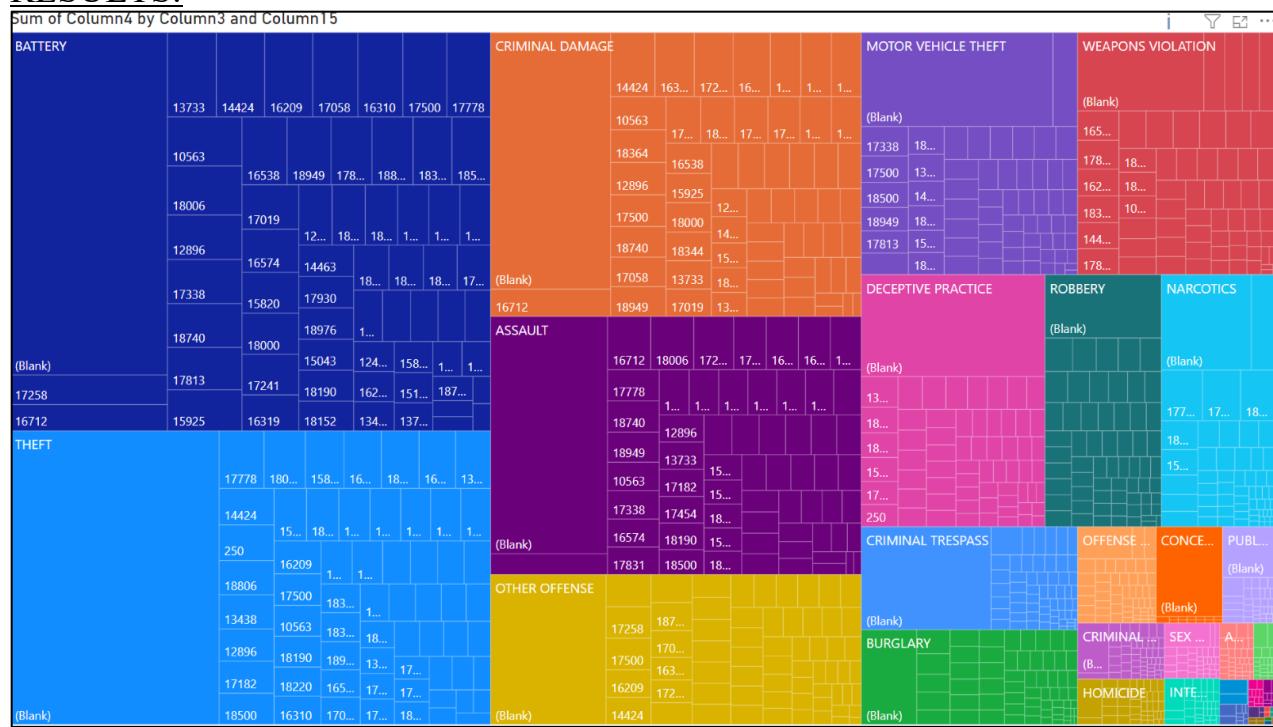
Category: Drag "primary_type" (column 3) to the Category field. This will create the main groupings or segments in the treemap based on different crime types.

6)Values: Drag "count" (or the metric representing the number of crimes)(column 4) to the Values field. This determines the size of each segment in the treemap based on the crime count.

7)Details : (column 15), you can drag median_income into details section

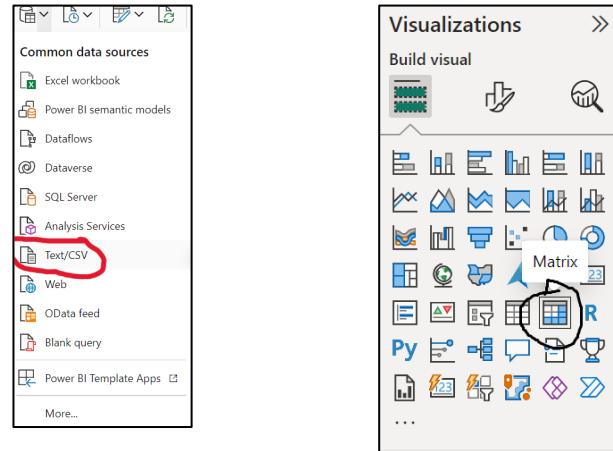


RESULTS:



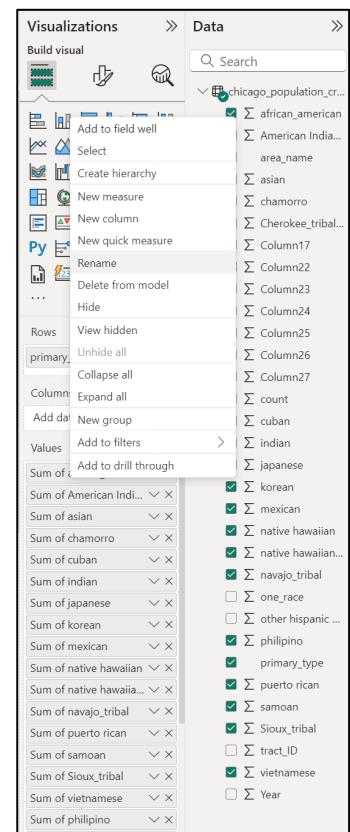
LOS ANGELES POPULATION HEAT MAP:

- STEP 1: Open a new PowerBI Blank Document
 STEP 2: in the GET DATA module under home click text/csv
 STEP 3: In the data preview window click load
 STEP 4: in the visualization, click the heat map symbol



STEP 5: Label columns correspondingly

COLUMN 1 = tract_id
 COLUMN 2 = year
 COLUMN 3 = crm_cd
 COLUMN 4 = count
 COLUMN 5 = one_race
 COLUMN 6 = black_or_african_american
 COLUMN 7 = american_indian_and_alaska_native
 COLUMN 8 = cherokee_tribal_grouping
 COLUMN 9 = chippewa_tribal_grouping
 COLUMN 10 = navajo_tribal_grouping
 COLUMN 11 = sioux_tribal_grouping
 COLUMN 12 = asian
 COLUMN 13 = indian
 COLUMN 14 = chinese
 COLUMN 15 = filipino
 COLUMN 16 = japanese
 COLUMN 17 = korean
 COLUMN 18 = vietnamese
 COLUMN 19 = other_asian
 COLUMN 20 = native_hawaiian_and_other_pacific_islander
 COLUMN 21 = chamorro
 COLUMN 22 = native_hawaiian
 COLUMN 23 = samoan COLUMN 24 = other_native_hawaiian_and_other_pacific_islander COLUMN 25 = race_alone_or_in_combination_with_one_or_more_other_races_total_population
 COLUMN 26 = hispanic_or_latino_total_population
 COLUMN 27 = hispanic_or_latino_total_population_of_any_race
 COLUMN 28 = mexican COLUMN 29 = puerto_rican
 COLUMN 30 = cuban
 COLUMN 31 = other_hispanic_or_latino
 COLUMN 32 = area_name COLUMN 33 = description



STEP 6:

In the "row" section put the
"description" here

STEP 7: Under the "values" put all the ethnicities column 6 - 31

Add data fields here		
Values	▼	×
Sum of black or africa...	▼	×
Sum of asian	▼	×
Sum of cherokee triba...	▼	×
Sum of chamorro	▼	×
Sum of chinese	▼	×
Sum of chippewa_trib...	▼	×
Sum of cuban	▼	×
Sum of indian	▼	×
Sum of filipino	▼	×
Sum of japanese	▼	×
Sum of korean	▼	×
Sum of native hawaiian	▼	×
Sum of mexican	▼	×
Sum of native_hawaiia...	▼	×
Sum of navajo_tribal	▼	×
Sum of samoan	▼	×
Sum of puerto_rican	▼	×

Result:

description	Sum of black or african american	Sum of asian	Sum of cherokee tribal grouping	Sum of chamorro	Sum of chinese	Sum of chij	Sum of ci
ARSON	281707	283988	600	951	49645		
ASSAULT WITH DEADLY WEAPON AGGRAVATED ASSAULT	680458	867343	1689	2565	151139		
ASSAULT WITH DEADLY WEAPON ON POLICE OFFICER	125056	86081	106	293	19734		
ATTEMPTED ROBBERY	431266	435576	818	1431	77369		
BATTERY - SIMPLE ASSAULT	700525	934257	1742	2664	164298		
BATTERY ON A FIREFIGHTER	26928	34030	87	23	7938		
BATTERY POLICE (SIMPLE)	238753	204791	277	1119	40236		
BATTERY WITH SEXUAL CONTACT	365527	437296	836	1256	81021		
BEASTIALITY CRIME AGAINST NATURE SEXUAL ASSLT WITH ANIM	1878	385	0	0	184		
BIGAMY	276	160	0	0	3		
BIKE - ATTEMPTED STOLEN	3481	3569	0	0	1381		
BIKE - STOLEN	300891	536331	916	1175	112158		
BLOCKING DOOR INDUCTION CENTER	630	430	9	0	31		
BOAT - STOLEN	10267	17986	23	13	2290		
BOMB SCARE	43466	68028	86	127	15709		
BRANDISH WEAPON	616639	724194	1385	2268	124432		
BRIBERY	915	790	0	0	85		
BUNCO ATTEMPT	73145	47777	65	303	8659		
BUNCO GRAND THEFT	427641	644046	1238	2049	114254		
BUNCO PETTY THEFT	252144	336545	698	1077	58854		
BURGLARY	690060	924309	1733	2572	163990		
BURGLARY ATTEMPTED	350017	475359	1180	1341	85387		
BURGLARY FROM VEHICLE	6905944	922079	1697	2680	163985		
BURGLARY FROM VEHICLE ATTEMPTED	109360	119114	266	606	22392		
CHILD ABANDONMENT	7022	4788	0	0	1451		
CHILD ABUSE (PHYSICAL) - AGGRAVATED ASSAULT	132106	76771	103	435	10612		
CHILD ABUSE (PHYSICAL) - SIMPLE ASSAULT	391879	339549	801	1448	53479		
CHILD ANNOYING (17YRS & UNDER)	138774	148011	267	651	25614		
CHILD NEGLECT (SEE 300 W.J.C.)	173222	162414	417	427	26636		
CHILD PORNOGRAPHY	44954	44403	94	86	9930		
CHILD STEALING	116493	61369	265	244	10836		
CONSPIRACY	2631	3727	0	19	271		
CRIMINAL HOMICIDE	245862	110467	264	477	15174		
Total	25048390	30043448	57980	90914	5334900		

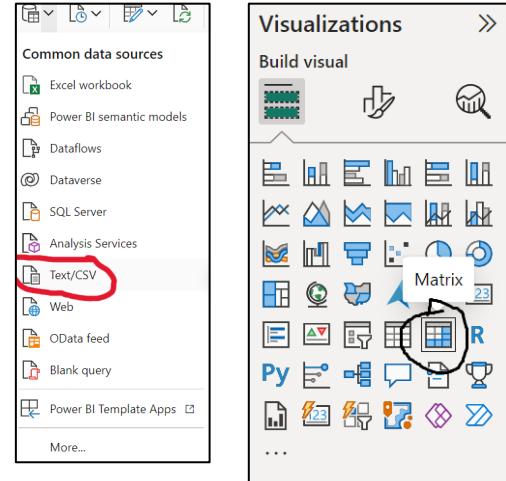
CHICAGO POPULATION HEAT MAP:

STEP 1: Open a new PowerBI Blank Document

STEP 2: in the GET DATA module under home click text/csv

STEP 3: In the data preview window click load

STEP 4: in the visualization, click the heat map symbol



STEP 5: Label columns correspondingly

COLUMN 1 = tract_id

COLUMN 2 = year

COLUMN 3 = primary_type

COLUMN 4 = count

COLUMN 5 = one_race

COLUMN 6 = black_or_african_american

COLUMN 7 = american_indian_and_alaska_native

COLUMN 8 = cherokee_tribal_grouping

COLUMN 9 = chippewa_tribal_grouping

COLUMN 10 = navajo_tribal_grouping

COLUMN 11 = sioux_tribal_grouping

COLUMN 12 = asian

COLUMN 13 = indian

COLUMN 14 = chinese

COLUMN 15 = filipino

COLUMN 16 = japanese

COLUMN 17 = korean

COLUMN 18 = vietnamese

COLUMN 19 = other_asian

COLUMN 20 = native_hawaiian_and_other_pacific_islander

COLUMN 21 = chamorro

COLUMN 22 = native_hawaiian

COLUMN 23 = samoan

COLUMN 24 = other_native_hawaiian_and_other_pacific_islander
COLUMN 25 = race_alone_or_in_combination_with_one_or_more_other_races_total_population

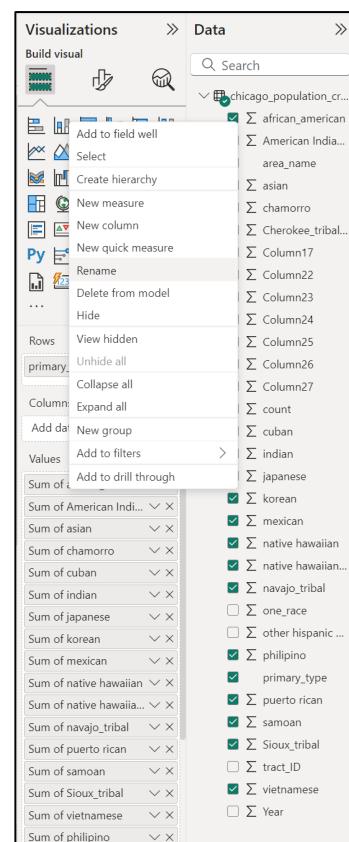
COLUMN 26 = hispanic_or_latino_total_population

COLUMN 27 = hispanic_or_latino_total_population_of_any_race
COLUMN 28 = mexican
COLUMN 29 = puerto_rican

COLUMN 30 = cuban

COLUMN 31 = other_hispanic_or_latino

COLUMN 32 = area_name



STEP 6:

In the "ROW" section put the "primary_type" here

STEP 7: Under the "values" put all the ethnicities column 6 - 31

Rows
primary_type
Columns
Add data fields here
Values
Sum of african_american
Sum of American_Indian
Sum of asian
Sum of chamorro
Sum of cuban
Sum of indian
Sum of japanese
Sum of korean
Sum of mexican
Sum of native_hawaiian
Sum of native_hawaiian
Sum of navajo_tribal
Sum of puerto_rican
Sum of samoan
Sum of Sioux_tribal
Sum of vietnamese
Sum of philipino

RESULT:

primary_type	m of mexican	Sum of native hawaiian	Sum of native hawaiian and other pacific islander	Sum of navajo_tribal	Sum of puerto rican	Sum of samoan	Sum of sioux_tribal	Sum of vietnamese	Sum of philipino
ASSAULT	585136	2337	21406	784	189303	507	1	1	1
BATTERY	587336	2352	21406	784	189462	522	1	1	1
BURGLARY	578996	2315	21305	784	188048	507	1	1	1
CRIMINAL DAMAGE	585771	2352	21406	784	189346	522	1	1	1
CRIMINAL TRESPASS	553740	2156	20149	783	174648	414	1	1	1
DECEPTIVE PRACTICE	584226	2352	21406	784	189780	522	1	1	1
MOTOR VEHICLE THEFT	585392	2342	21356	784	188615	522	1	1	1
NARCOTICS	471499	1747	16825	733	158422	398	1	1	1
OTHER OFFENSE	583914	2352	21393	784	188518	522	1	1	1
ROBBERY	572061	2214	20590	782	181991	490	1	1	1
SEX OFFENSE	411067	1573	15270	358	130906	429	1	1	1
THEFT	587336	2352	21406	784	189972	522	1	1	1
WEAPONS VIOLATION	555710	1929	17945	765	163831	465	1	1	1
ARSON	232345	733	7770	642	71898	197	1	1	1
OFFENSE INVOLVING CHILDREN	477165	1950	14917	679	145627	479	1	1	1
PUBLIC PEACE VIOLATION	250304	1192	12005	411	97770	355	1	1	1
CRIMINAL SEXUAL ASSAULT	421890	1686	16122	542	132242	399	1	1	1
HOMICIDE	247689	677	6405	417	64710	126	1	1	1
LIQUOR LAW VIOLATION	114331	459	3499	0	37874	188	1	1	1
STALKING	192825	640	5910	236	56365	89	1	1	1
INTIMIDATION	144546	415	5092	261	37465	121	1	1	1
INTERFERENCE WITH PUBLIC OFFICER	151797	630	4390	305	56783	154	1	1	1
OBSCENITY	39982	141	2545	191	11293	0	1	1	1
KIDNAPPING	81526	245	2541	59	23396	46	1	1	1
CONCEALED CARRY LICENSE VIOLATION	49060	234	718	159	18618	78	1	1	1
CRIM SEXUAL ASSAULT	589	147	2003	191	6898	26	1	1	1
GAMBLING	3948	52	494	0	2148	9	1	1	1
HUMAN TRAFFICKING	1758	0	411	0	1819	0	1	1	1
NON-CRIMINAL	522	0	63	0	118	0	1	1	1
OTHER NARCOTIC VIOLATION	7139	20	525	24	2322	11	1	1	1
PROSTITUTION	20637	164	239	0	7903	34	1	1	1
PUBLIC INDECENCY	1320	35	105	0	1130	1	1	1	1
DILITIGATION	0	0	0	0	0	0	1	1	1
Total	9681557	37793	347617	13810	3099221	8655	37.		

References

1. URL of Data Source: <https://data.census.gov/>
2. URL of Data Source: https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-Present/ijzp-q8t2/about_data
3. URL of Data Source: https://data.lacity.org/Public-Safety/Crime-Data-from-2020-to-Present/2nrs-mtv8/about_data
4. URL of Data Source: <https://www.esri.com/arcgis-blog/products/product/data-management/new-spatial-aggregation-tutorial-for-gis-tools-for-hadoop/?rmedium=redirect&rsource=blogs.esri.com/esri/arcgis/2015/03/25/new-spatial-aggregation-tutorial-for-gis-tools-for-hadoop>
5. GitHub: <https://github.com/xuemengtang/5200group2>