# Los Angeles and New York crime dataset analysis

**Course: BIG DATA TOOLS & TECHNO. II**
**Prof. Nizar Ali**

**Group 12 members:**

**Jose Anicio Pereira Junior - ID: 101506160**
**Himanshu Kakkar - ID: 101510716**
**Pujan Bhatt - ID: 101527283**

# Table of contents

- [Project overview](#)
- [Objectives](#)
- [Dataset and sources](#)
- [Data flow chart](#)
- [Data Preprocessing and transformation](#)
- [Results and graphs](#)
- [Conclusion](#)

# Project overview

This project focuses on analyzing and comparing crime rates in Los Angeles and New York City using extensive datasets provided by the respective police departments.

The aim is to integrate and process these datasets using various big data tools and technologies to uncover patterns and insights into crime dynamics in these two major cities.

# Objectives

The primary objective of this project is to compare crime rates between Los Angeles and New York City by analyzing large datasets.

This analysis aims to identify patterns, trends, and insights into the crime dynamics in these two major cities.

# Dataset and sources

**LYPD DATA**

Arrest Data from 2010 to 2019 -: https://catalog.data.gov/dataset/arrest-data-from-2010-to-2019

Crime Data from 2020 to Present -:https://catalog.data.gov/dataset/crime-data-from-2020-to-present

LAPD Calls for Service 2019-: https://catalog.data.gov/dataset/lapd-calls-for-service-2019
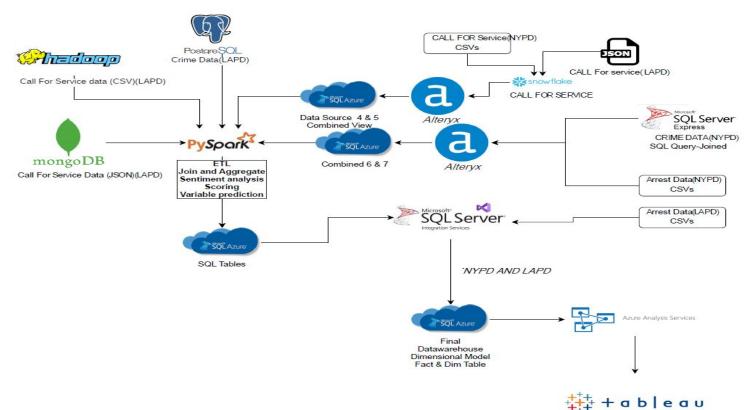
**NEW YORK DATA**

Arrest Data -: https://catalog.data.gov/dataset/nypd-arrests-data-historic

Call for Service Data -: https://catalog.data.gov/dataset/nypd-calls-for-service-historic

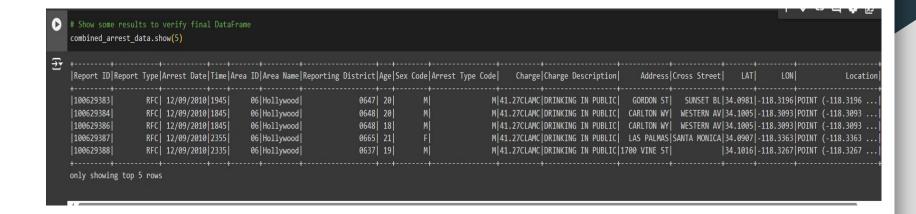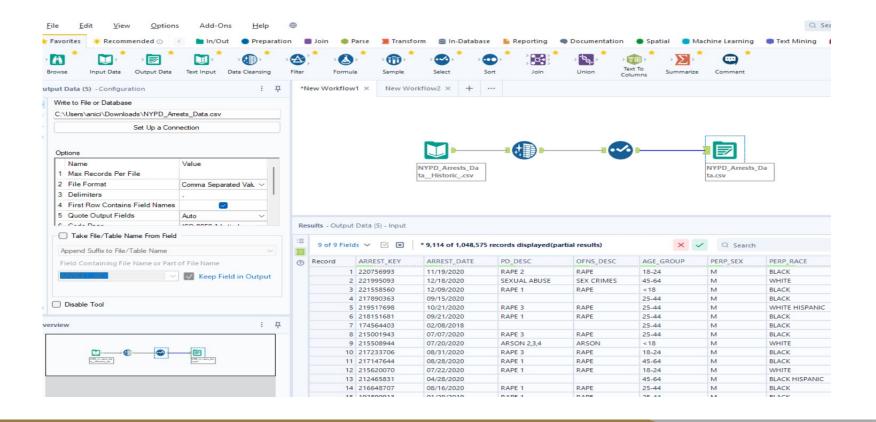Crime Data -: https://catalog.data.gov/dataset/nypd-complaint-data-historic

# Data flow chart

# Data Preprocessing and transformation

```python
# Function to clean address
def clean_address(address):
    if address:
        return ' '.join(address.split())
    return address

clean_address_udf = udf(clean_address, StringType())
```

```python
# Apply the cleaning function to the Address and Cross Street columns
combined_arrest_data = combined_arrest_data.withColumn('Address', clean_address_udf(col('Address')))
combined_arrest_data = combined_arrest_data.withColumn('Cross Street', clean_address_udf(col('Cross Street')))
```

```python
# Drop specified columns
columns_to_drop = [
    'Descent Code', 'Charge Group Code', 'Charge Group Description',
    'Disposition Description', 'Booking Date', 'Booking Time',
    'Booking Location', 'Booking Location Code'
]
combined_arrest_data = combined_arrest_data.drop(*columns_to_drop)
```

```python
# Fill missing values in 'Charge Description, cross street'
combined_arrest_data = combined_arrest_data.withColumn('Charge Description', when(col('Charge Description').isNull(), 'Unknown').otherwise(col('Charge Description')))
```

```python
combined_arrest_data = combined_arrest_data.withColumn('Cross Street', when(col('Cross Street').isNull(), 'NAN').otherwise(col('Cross Street')))
# Fill missing values with an empty string
combined_arrest_data = combined_arrest_data.na.fill('NAN')
```

```python
# Show some results to verify final DataFrame
combined_arrest_data.show(5)
```
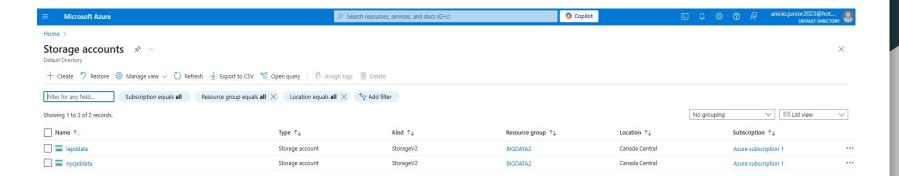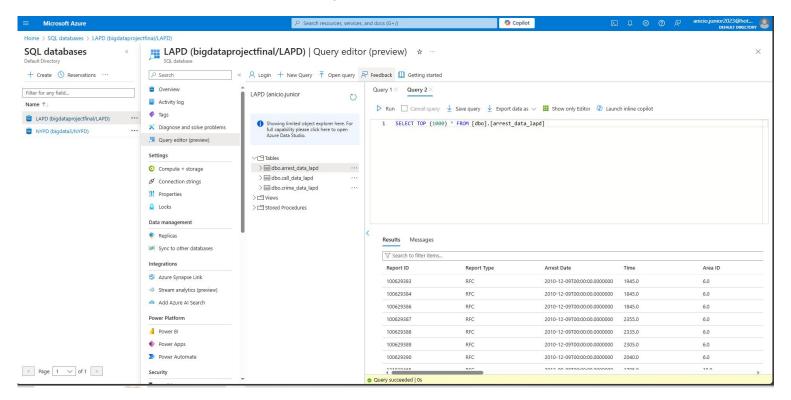
# Data Preprocessing and transformation
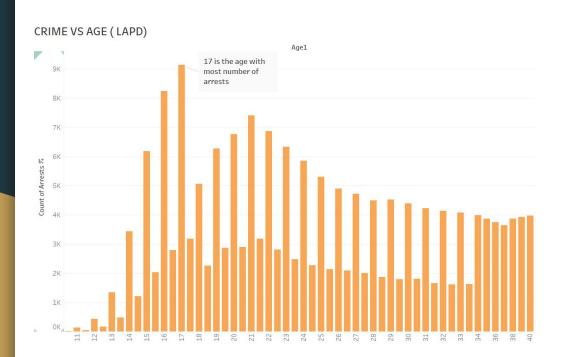
# Data Preprocessing and transformation

# Azure connectivity

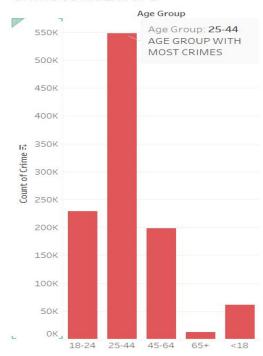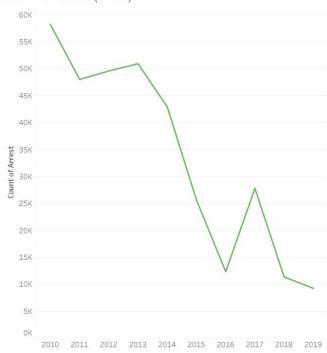# Azure connectivity

# Results and graphs



CRIME VS AGE ( LAPD)

17 is the age with most number of arrests



Crime vs AGE NYPD

Age Group: **25-44** AGE GROUP WITH MOST CRIMES

# Results and graphs



Arrest VS YEAR (LAPD)



ARREST VS YEAR (NYPD)

# Results and graphs

Most Calls For Service (LAPD)

| | | | | | |
|---|---|---|---|---|---|
| CODE 6 | TRAFFIC STOP | MAN | GRP | CODE 30 RINGER | SUSP NOW |
| | TRESPASS SUSP | RADIO | | POSS SUSP | SUSP J/L |
| | INVEST | J/O | | | |
| | SUSP | H & R MISD | CPI | | |
| | | | I/P | | |
| | AMB | | OTHER | | |
| | DOM VIOL | | FIGHT | | |
| | | BUSN | MALE | | |
| | PARTY | WMN | | | |

Most CALLS for Service(NYPD)

| | | | | | |
|---|---|---|---|---|---|
| VISIBILITY PATROL: DIRECTED | TRAIN RUN/MOBILE ORDER MAINTENANCE SWEEP | VISIBILITY PATROL: FAMILY/HOME VISIT | | ALARMS: | DISPUTE: INSIDE | TRAIN ORDER |
| | TRANSIT PATROL/INSPECTION BY NON-TRANSIT BUREAU PERSONNEL | VISIBILITY PATROL: INTERIOR | | | |
| STATION INSPECTION BY TRANSIT BUREAU PERSONNEL | INVESTIGATE/POSSIBLE CRIME: SUSP VEHICLE/OUTSIDE | VEHICLE ACCIDENT: SPECIAL | SEE | | |
| | | TRAFFIC SAFETY | VERIFY AMB | | |
| | INVESTIGATE/POSSIBLE CRIME: CALLS FOR HELP/INSIDE | OTHER CRIMES (IN PROGRESS): | | | |
| SEE COMPLAINANT: OTHER/INSIDE | | DISORDERLY: | | | |
| | AMBULANCE CASE: EDP/INSIDE | AMBULANCE CASE: | | | |

# Results and graphs

NUMBER OF CALLS FOR CRIME VS PERSON ( SEX/ RACE)

# Results and graphs



MOST CRIMES (LAPD)

SIT/LIE/SLEEP SIDEWALK OR STREET
TRUANCY
ILLEGAL
DRINKING IN PUBLIC
CURFEW - JUV ONLY
LOS ANGELES MUNICIPAL
OPEN ALCOHOLIC BEV IN PUBLIC PARK/PLACE
POSS ALCH BEV ON

MOST CRIMES(NYPD)

DANGEROUSDRUGS
FELONYASSAULT
GRANDLARCENY
MISCELLANEOUSPENALLAW
BURGLARY
FORGERY
ASSAULT3&RELATEDOFFENSES
VEHICLEANDTRAFFICLAWS
ROBBERY
RAPE
SEXCRIMES
PETITLARCENY
DANGEROUSWEAPONS

# Conclusion

**Key Findings:**

- **Crime vs. Age:**
    - **LAPD:** Peak arrests at age 17, smaller peak around ages 23-24.
    - **NYPD:** Highest crime count in the 25-44 age group.
- **Arrests vs. Year:**
    - **LAPD:** Overall declining trend from 2010 to 2019.
    - **NYPD:** Significant drop around 2012, followed by fluctuations.
- **Most Calls for Service:**
    - **LAPD:** 'Code 6' is the most frequent call.
    - **NYPD:** 'Visibility Patrol Directed' is the most common call.

**Implications:**

- Differences in age-related crime trends suggest targeted interventions.
- Yearly arrest trends indicate varying impacts of policies and socio-economic factors.
- Most common service calls reflect different policing priorities and strategies.

**Conclusion:**

- Understanding these patterns can help optimize law enforcement strategies and improve crime prevention efforts.
- Insights are valuable for developing targeted interventions in both cities.