

# Cloud Computing

The screenshot displays a Databricks notebook titled "Project#4: Yasaswi\_CC\_Cluster" in a web browser. The interface includes a sidebar with navigation icons, a top toolbar with "Run all", "Terminated", and "Publish" buttons, and a main workspace area.

**Cell 4: Getting Started**

Run the following cell to configure our "classroom."

```
1 %run ../Includes/Classroom-Setup
```

Command took 3.21 seconds -- by a user at 16/4/2023, 7:44:43 pm on unknown cluster

Initialized classroom variables & functions...

Datasets are already mounted to /mnt/training from wasbs://training@dbtraineastus.blob.core.windows.net/

Imported Test Library..

Created user specific database

Using the database kurrayi\_mail\_uc\_edu\_db

All done!

**Cell 5:**

```
1 print("username: " + username)
2 print("userhome: " + userhome)
```

username: kurrayi@mail.uc.edu  
userhome: dbfs:/user/kurrayi@mail.uc.edu

Command took 0.18 seconds -- by a user at 16/4/2023, 7:44:43 pm on unknown cluster

**Cell 6:**

```
2 crimeDataLosAngelesDF = spark.read.parquet('/mnt/training/crime-data-2016/Crime-Data-Los-Angeles-2016.parquet')
3 print("Total count of the crimes occurred in Los Angeles : ",crimeDataLosAngelesDF.count())
4
5 crimeDataPhiladelphiaDF = spark.read.parquet("/mnt/training/crime-data-2016/Crime-Data-Philadelphia-2016.parquet")
6 print("Total count of the crimes occurred in Philadelphia : ",crimeDataPhiladelphiaDF.count())
7
8 crimeDataDallasDF = spark.read.parquet("/mnt/training/crime-data-2016/Crime-Data-Dallas-2016.parquet")
9 print("Total count of the crimes occurred in Dallas : ",crimeDataDallasDF.count())
10
11 crimeDataNewYorkDF = spark.read.parquet("/mnt/training/crime-data-2016/Crime-Data-New-York-2016.parquet")
12 print("Total count of the crimes occurred in New York : ",crimeDataNewYorkDF.count())
13
14 crimeDataChicagoDF = spark.read.parquet("/mnt/training/crime-data-2016/Crime-Data-Chicago-2016.parquet")
15 print("Total count of the crimes occurred in Chicago : ",crimeDataChicagoDF.count())
16
17 crimeDataBostonDF = spark.read.parquet("/mnt/training/crime-data-2016/Crime-Data-Boston-2016.parquet")
18 print("Total count of the crimes occurred in Boston : ",crimeDataBostonDF.count())
```

(18) Spark Jobs

- crimeDataLosAngelesDF: pyspark.sql.dataframe.DataFrame = [id: string, dateReported: timestamp ... 17 more fields]
- crimeDataPhiladelphiaDF: pyspark.sql.dataframe.DataFrame = [district: integer, dispatch\_date\_time: timestamp ... 12 more fields]
- crimeDataDallasDF: pyspark.sql.dataframe.DataFrame = [incidentNumberWithYear: string, incidentNumberWithoutYear: integer ... 78 more fields]
- crimeDataNewYorkDF: pyspark.sql.dataframe.DataFrame = [complaintNumber: integer, keyCode: integer ... 16 more fields]
- crimeDataChicagoDF: pyspark.sql.dataframe.DataFrame = [id: integer, caseNumber: string ... 20 more fields]
- crimeDataBostonDF: pyspark.sql.dataframe.DataFrame = [INCIDENT\_NUMBER: string, OFFENSE\_CODE: integer ... 15 more fields]

Total count of the crimes occurred in Los Angeles : 217945  
Total count of the crimes occurred in Philadelphia : 168664  
Total count of the crimes occurred in Dallas : 99642  
Total count of the crimes occurred in New York : 468241  
Total count of the crimes occurred in Chicago : 267872  
Total count of the crimes occurred in Boston : 218610

Command took 31.40 seconds -- by kurrayi@mail.uc.edu at 13/4/2023, 9:45:07 pm on Yasaswi\_CC\_Cluster

University of Cincinnati Project#4: Yasaswi\_CC\_Cluster Python

File Edit View Run Help Last edit was yesterday Give feedback

Run all Terminated Publish

(24) Spark Jobs

- crimeDataLosAngelesDF: pyspark.sql.dataframe.DataFrame = [id: string, dateReported: timestamp ... 17 more fields]
- crimeDataPhiladelphiaDF: pyspark.sql.dataframe.DataFrame = [district: integer, dispatch\_date\_time: timestamp ... 12 more fields]
- crimeDataDallasDF: pyspark.sql.dataframe.DataFrame = [incidentNumberWithYear: string, incidentNumberWithoutYear: integer ... 78 more fields]
- crimeDataNewYorkDF: pyspark.sql.dataframe.DataFrame = [complaintNumber: integer, keyCode: integer ... 16 more fields]
- crimeDataChicagoDF: pyspark.sql.dataframe.DataFrame = [id: integer, caseNumber: string ... 20 more fields]
- crimeDataBostonDF: pyspark.sql.dataframe.DataFrame = [INCIDENT\_NUMBER: string, OFFENSE\_CODE: integer ... 15 more fields]
- x: pyspark.sql.dataframe.DataFrame = [count(DISTINCT OFFENSE\_CODE\_GROUP): long]

Total count of different types of the crimes occurred in Los Angeles :

```
|count(DISTINCT crimeCodeDescription)|
|-----|
|124|
```

Total count of different types of the crimes occurred in Philadelphia :

```
|count(DISTINCT ucr_general)|
|-----|
|26|
```

Total count of different types of the crimes occurred in Dallas :

```
|count(DISTINCT typeOfIncident)|
|-----|
|378|
```

Command took 35.45 seconds -- by kurrayi@mail.uc.edu at 13/4/2023, 9:46:59 pm on Yasaswi\_CC\_Cluster

Ced 12

University of Cincinnati Project#4: Yasaswi\_CC\_Cluster Python

File Edit View Run Help Last edit was yesterday Give feedback

Run all Terminated Publish

```
crimeDataPhiladelphiaDF robberyPhiladelphiaDF ucr_general_description
crimeDataDallasDF robberyDallasDF typeOfIncident
```

Ced 13

```
1 from pyspark.sql.functions import col, lower
2 LA = crimeDataLosAngelesDF.select(lower(col("crimeCodeDescription")))
3 robberyLosAngelesDF = LA.filter("lower(crimeCodeDescription) like 'robbery%'").count()
4 print("Total robbery count occurred in Los Angeles :", robberyLosAngelesDF)
5
6 P = crimeDataPhiladelphiaDF.select(lower(col("ucr_general_description")))
7 robberyPhiladelphiaDF = P.filter("lower(ucr_general_description) like 'robbery%'").count()
8 print("Total robbery count occurred in Philadelphia:", robberyPhiladelphiaDF)
9
10 D = crimeDataDallasDF.select(lower(col("typeOfIncident")))
11 robberyDallasDF = D.filter("lower(typeOfIncident) like 'robbery%'").count()
12 print("Total robbery count occurred in Dallas:", robberyDallasDF)
```

(6) Spark Jobs

- LA: pyspark.sql.dataframe.DataFrame = [lower(crimeCodeDescription): string]
- P: pyspark.sql.dataframe.DataFrame = [lower(ucr\_general\_description): string]
- D: pyspark.sql.dataframe.DataFrame = [lower(typeOfIncident): string]

Total robbery count occurred in Los Angeles : 18275  
Total robbery count occurred in Philadelphia: 6149  
Total robbery count occurred in Dallas: 6852

Command took 21.48 seconds -- by kurrayi@mail.uc.edu at 13/4/2023, 9:47:25 pm on Yasaswi\_CC\_Cluster

Ced 14

University of Cincinnati Project#4: Ysaswi\_CC\_Cluster Python

File Edit View Run Help Last edit was yesterday Give feedback

Run all Terminated Publish

```
16 print("\n(Philadelphia) Highest is", robberiesByMonthPhiladelphiaDF.head(1))
17 print("(Philadelphia) Lowest is", robberiesByMonthPhiladelphiaDF.tail(1))
18 robberyDallasDF = crimeDataDallasDF.withColumn("monthOccurred", split(col("callReceivedDateTime"), "/").getItem(0)).withColumn("crime", lower(col("typeOfIncident")))
19 robberiesByMonthDallasDF = robberyDallasDF.select("monthOccurred", "crime")
20 robberiesByMonthDallasDF = robberiesByMonthDallasDF.filter("crime like 'robbery%")
21 robberiesByMonthDallasDF = robberiesByMonthDallasDF.groupby("monthOccurred").count().sort(desc('count'))
22 print("\n(Dallas) Highest is", robberiesByMonthDallasDF.head(1))
23 print("(Dallas) Lowest is", robberiesByMonthDallasDF.tail(1))
```

(18) Spark Jobs

- robberyLosAngelesDF: pyspark.sql.dataframe.DataFrame = [id: string, dateReported: timestamp ... 17 more fields]
- df1: pyspark.sql.dataframe.DataFrame = [id: string, dateReported: timestamp ... 17 more fields]
- robberyLosAngelesDF1: pyspark.sql.dataframe.DataFrame = [timeOccurred: integer, crimeCodeDescription: string]
- robberiesByMonthLosAngelesDF: pyspark.sql.dataframe.DataFrame = [timeOccurred: integer, count: long]
- robberiesByMonthDallasDF: pyspark.sql.dataframe.DataFrame = [district: integer, dispatch\_date\_time: timestamp ... 12 more fields]
- robberiesByMonthDallasDF1: pyspark.sql.dataframe.DataFrame = [dispatch\_date\_time: integer, ucr\_general\_description: string]
- robberiesByMonthDallasDF2: pyspark.sql.dataframe.DataFrame = [dispatch\_date\_time: integer, count: long]
- robberiesByMonthDallasDF3: pyspark.sql.dataframe.DataFrame = [dispatch\_date\_time: integer, count: long]
- robberyDallasDF: pyspark.sql.dataframe.DataFrame = [incidentNumberWithYear: string, incidentNumberWithoutYear: integer ... 80 more fields]
- robberiesByMonthDallasDF: pyspark.sql.dataframe.DataFrame = [monthOccurred: string, count: long]

(Los Angeles) Highest is [Row(timeOccurred=12, count=853)]  
(Los Angeles) Lowest is [Row(timeOccurred=2, count=675)]

(Philadelphia) Highest is [Row(dispatch\_date\_time=10, count=572)]  
(Philadelphia) Lowest is [Row(dispatch\_date\_time=2, count=416)]

(Dallas) Highest is [Row(monthOccurred='01', count=743)]  
(Dallas) Lowest is [Row(monthOccurred='03', count=410)]

Command took 55.40 seconds -- by kurrayi@mail.uc.edu at 13/4/2023, 9:47:33 pm on Ysaswi\_CC\_Cluster

Ced 18

University of Cincinnati Project#4: Ysaswi\_CC\_Cluster Python

File Edit View Run Help Last edit was yesterday Give feedback

Run all Terminated Publish

**Hint:** You may want to apply the `union()` method in this example to combine the three datasets.

Ced 19

**Hint:** Combined 3 Cities

Ced 20

```
1 # TODO
2
3 from pyspark.sql.functions import lit
4 robberiesByMonthLosAngelesDF = robberiesByMonthLosAngelesDF.withColumn("city", lit("Los Angeles"))
5 robberiesByMonthPhiladelphiaDF = robberiesByMonthPhiladelphiaDF.withColumn("city", lit("Philadelphia"))
6 robberiesByMonthDallasDF = robberiesByMonthDallasDF.withColumn("city", lit("Dallas"))
7 combinedRobberiesByMonthDF = robberiesByMonthDallasDF.union(robberiesByMonthLosAngelesDF).union(robberiesByMonthPhiladelphiaDF)
8 high = combinedRobberiesByMonthDF.sort(desc('count')).head(1)
9 low = combinedRobberiesByMonthDF.sort(desc('count')).tail(1)
10 print("Highest is at : ", high)
11 print("Lowest is at : ", low)
```

(10) Spark Jobs

- robberiesByMonthLosAngelesDF: pyspark.sql.dataframe.DataFrame = [timeOccurred: integer, count: long ... 1 more field]
- robberiesByMonthPhiladelphiaDF: pyspark.sql.dataframe.DataFrame = [dispatch\_date\_time: integer, count: long ... 1 more field]
- robberiesByMonthDallasDF: pyspark.sql.dataframe.DataFrame = [monthOccurred: string, count: long ... 1 more field]
- combinedRobberiesByMonthDF: pyspark.sql.dataframe.DataFrame = [monthOccurred: string, count: long ... 1 more field]

Highest is at : [Row(monthOccurred='12', count=853, city='Los Angeles')]  
Lowest is at : [Row(monthOccurred='03', count=410, city='Dallas')]

Command took 48.35 seconds -- by kurrayi@mail.uc.edu at 13/4/2023, 9:47:41 pm on Ysaswi\_CC\_Cluster

Ced 21

Project#4: Yasaswi\_CC\_Cluster Python

```

9 Philadelphiadata = crimeDataPhiladelphiaDF.withColumn("monthOccurred",split(col("dispatch_date"),"-").getItem(1)).withColumn("crime",lower(col("ucr_general_description")))
10 robberiesByMonthPhiladelphiaDF = Philadelphiadata.select("monthOccurred","crime")
11 robberiesByMonthPhiladelphiaDF = robberiesByMonthPhiladelphiaDF.filter("crime like '%robbery%'")
12
13 Dallasdata = crimeDataDallasDF.withColumn("monthOccurred",split(col("callReceivedDate"),"/").getItem(0)).withColumn("crime",lower(col("typeOfIncident")))
14 robberiesByMonthDallasDF = Dallasdata.select("monthOccurred","crime")
15 robberiesByMonthDallasDF = robberiesByMonthDallasDF.filter("crime like '%robbery%'")
16
17 robbery_per_month_dallas_city = robberiesByMonthDallasDF.groupBy("monthOccurred").count().withColumn("city",lit("Dallas"))
18 robbery_per_month_losangeles_city = robberiesByMonthLosAngelesDF.groupBy("monthOccurred").count().withColumn("city",lit("LosAngeles"))
19 robbery_per_month_philadelphia_city = robberiesByMonthPhiladelphiaDF.groupBy("monthOccurred").count().withColumn("city",lit("Philadelphia"))
20 combinedRobberiesByMonthDF = robbery_per_month_dallas_city.union(robbery_per_month_losangeles_city).union(robbery_per_month_philadelphia_city).withColumnRenamed("count","robbery").sort("monthOccurred",ascending=True)
21
22 display(combinedRobberiesByMonthDF)

```

(4) Spark Jobs

- LosAngelesData: pyspark.sql.dataframe.DataFrame = [id: string, dateReported: timestamp ... 19 more fields]
- robberiesByMonthLosAngelesDF: pyspark.sql.dataframe.DataFrame = [monthOccurred: string, crime: string]
- Philadelphiadata: pyspark.sql.dataframe.DataFrame = [district: integer, dispatch\_date: timestamp ... 14 more fields]
- robberiesByMonthPhiladelphiaDF: pyspark.sql.dataframe.DataFrame = [monthOccurred: string, crime: string]
- Dallasdata: pyspark.sql.dataframe.DataFrame = [incidentNumberWithYear: string, incidentNumberWithoutYear: integer ... 80 more fields]
- robberiesByMonthDallasDF: pyspark.sql.dataframe.DataFrame = [monthOccurred: string, crime: string]
- robbery\_per\_month\_dallas\_city: pyspark.sql.dataframe.DataFrame = [monthOccurred: string, count: long ... 1 more field]
- robbery\_per\_month\_losangeles\_city: pyspark.sql.dataframe.DataFrame = [monthOccurred: string, count: long ... 1 more field]
- robbery\_per\_month\_philadelphia\_city: pyspark.sql.dataframe.DataFrame = [monthOccurred: string, count: long ... 1 more field]
- combinedRobberiesByMonthDF: pyspark.sql.dataframe.DataFrame = [monthOccurred: string, robbery: long ... 1 more field]

Project#4: Yasaswi\_CC\_Cluster Python

combinedRobberiesByMonthDF: pyspark.sql.dataframe.DataFrame = [monthOccurred: string, robbery: long ... 1 more field]

	monthOccurred	robbery	city
1	01	743	Dallas
2	01	835	LosAngeles
3	01	524	Philadelphia
4	02	757	LosAngeles
5	02	439	Dallas
6	02	412	Philadelphia

36 rows | 24.35 seconds runtime

Command took 24.35 seconds -- by kurrayi@mail.uc.edu at 13/4/2023, 9:47:51 pm on Yasaswi\_CC\_Cluster

**Question 7: Find the "per capita robbery rates" using ther Hint below, and plot graph as above for the per capita robbery rates per month for each of the three cities, producing one Graph using the contents of robberyRatesByCityDF :**

- Los Angeles
- Philadelphia
- Dallas

While the above graph is interesting, it's flawed: it's comparing the raw numbers of robberies, not the per capita robbery rates.

University of Cincinnati | Project#4: Ysaswi\_CC\_Cluster | Python | Run all | Terminated | Publish

```
8 population_LosAngeles = cityDataDF.filter("city = 'Los Angeles'").select("estPopulation2016").collect()[0]["estPopulation2016"]
9 population_Philadelphia = cityDataDF.filter("city = 'Philadelphia'").select("estPopulation2016").collect()[0]["estPopulation2016"]
10 population_Dallas = cityDataDF.filter("city = 'Dallas'").select("estPopulation2016").collect()[0]["estPopulation2016"]
11 robbery_rates_for_LA = combinedRobberiesByMonthDF.withColumn("robberyRate", (combinedRobberiesByMonthDF.robbery)/population_LosAngeles).where("city = 'Los Angeles'")
12 robbery_rates_for_Philadelphia = combinedRobberiesByMonthDF.withColumn("robberyRate", (combinedRobberiesByMonthDF.robbery)/population_Philadelphia).where("city = 'Philadelphia'")
13 robbery_rates_for_Dallas = combinedRobberiesByMonthDF.withColumn("robberyRate", (combinedRobberiesByMonthDF.robbery)/population_Dallas).where("city = 'Dallas'")
14 robberyRatesByCityDF = robbery_rates_for_LA.union(robbery_rates_for_Philadelphia).union(robbery_rates_for_Dallas).sort("monthOccurred", ascending=True)
15 display(robberyRatesByCityDF)
```

Table

	monthOccurred	robbery	city	robberyRate
1	01	743	Dallas	0.0005637632983263893
2	01	835	LosAngeles	0.0002099930538824572
3	01	524	Philadelphia	0.0003342109559964079
4	02	757	LosAngeles	0.00019037693627427558
5	02	439	Dallas	0.00033309836872851266
6	02	412	Philadelphia	0.00026277655318801535

36 rows | 25.61 seconds runtime | Refreshed 3 days ago

Command took 25.61 seconds -- by kurrayi@mail.uc.edu at 13/4/2023, 9:47:57 pm on Ysaswi\_CC\_Cluster

University of Cincinnati | Project#4: Ysaswi\_CC\_Cluster | Python | Run all | Terminated | Publish

```
26 combinedHomicidesByMonthDF = reduce(DataFrame.unionAll, dfs)
27 combinedHomicidesByMonthDF1 = combinedHomicidesByMonthDF.sort(desc('count'))
28 print(combinedHomicidesByMonthDF1.head(1))
29 print(combinedHomicidesByMonthDF1.tail(1))
```

Row(month=8, count=36, city='New York')  
Row(month=3, count=4, city='Boston')

Command took 10.13 seconds -- by kurrayi@mail.uc.edu at 13/4/2023, 9:48:04 pm on Ysaswi\_CC\_Cluster

Cmd 27

**Question 9: Find the "per capita HOMICIDE rates" using the Hint in Qn 7, and plot graph as above for the per capita HOMICIDE rates per month for each of the two cities, producing one Graph using the content of HOMICIDE rates per month for each of the two cities, producing one Graph using**



University of Cincinnati

Project#4: Ysaswi\_CC\_Cluster

https://community.cloud.databricks.com/?o=2784418577973442#notebook/2083315109028303/command/2083315109028330

Run allTerminatedPublish

File Edit View Run Help Last edit was yesterday Give feedback

ced 28

```
1 from pyspark.sql import functions as F
2 from pyspark.sql.functions import month, desc
3
4 cityDataDF= spark.read.parquet("dbfs:/mnt/training/City-Data.parquet")
5 HOMICIDERatesByCityDF= combinedHomicidesByMonthDF.join(cityDataDF, on=['city'], how='inner').withColumn('HomicideRate', F.col('count')/F.col('estPopulation2016'))
6
7 HOMICIDERatesByCityDF=HOMICIDERatesByCityDF.select(F.col("month"),F.col("HomicideRate"),F.col("city")).orderBy(F.col("month"))
8 display(HOMICIDERatesByCityDF)
```

(5) Spark Jobs

cityDataDF: pyspark.sql.dataframe.DataFrame = [rankin2016: integer, state: string ... 4 more fields]

HOMICIDERatesByCityDF: pyspark.sql.dataframe.DataFrame = [month: integer, HomicideRate: double ... 1 more field]

Table +

	month	HomicideRate	city
1	1	0.0000025768145488823477	New York
2	1	0.000010398345771735514	Boston
3	2	0.0000019911748786818143	New York
4	2	0.00000594191186956315	Boston
5	3	0.000002928198351002668	New York
6	3	0.00000594191186956315	Boston

24 rows | 6.04 seconds runtime

Refreshed 3 days ago

Command took 6.04 seconds -- by kurray1@mail.uc.edu at 13/4/2023, 9:48:09 pm on Ysaswi\_CC\_Cluster

University of Cincinnati

Project#4: Ysaswi\_CC\_Cluster

https://community.cloud.databricks.com/?o=2784418577973442#notebook/2083315109028303/command/2083315109028330

Run allTerminatedPublish

File Edit View Run Help Last edit was yesterday Give feedback

ced 30

1 The models, which comprise machine learning, artificial neural networks, support vector machines, and fuzzy logic, should make it possible to forecast time series models. The two main goals of time series analysis are:

2 One involves acknowledging the possibility that the event covered by the insight collection will occur, and the other comprises preparing for upcoming evaluations of the time arrangement variable.

3 With historical data from crime-data-2016, it is possible to project future values for the monthly Robbery Count rate for Log Angeles by visualizing the time series, normalizing the time series, finding the ideal model parameters, fitting the model, and creating forecasts.

ced 31

## References

The crime data used in this notebook comes from the following locations:

City	Original Data
Boston	<a href="https://data.boston.gov/group/public-safety">https://data.boston.gov/group/public-safety</a>
Chicago	<a href="https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-present/ijzp-q8t2">https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-present/ijzp-q8t2</a>
Dallas	<a href="https://www.dallasopendata.com/Public-Safety/Police-Incidents/tbnj-w5hb/data">https://www.dallasopendata.com/Public-Safety/Police-Incidents/tbnj-w5hb/data</a>
Los Angeles	<a href="https://data.lacity.org/A-Safe-City/Crime-Data-From-2010-to-Present/y8tr-7khq">https://data.lacity.org/A-Safe-City/Crime-Data-From-2010-to-Present/y8tr-7khq</a>
New Orleans	<a href="https://data.nola.gov/Public-Safety-and-Preparedness/Electronic-Police-Report-2016/4gc2-25he/data">https://data.nola.gov/Public-Safety-and-Preparedness/Electronic-Police-Report-2016/4gc2-25he/data</a>
New York	<a href="https://data.cityofnewyork.us/Public-Safety/NYPD-Complaint-Data-Historic/qgea-i56i">https://data.cityofnewyork.us/Public-Safety/NYPD-Complaint-Data-Historic/qgea-i56i</a>
Philadelphia	<a href="https://www.opendataphilly.org/dataset/crime-incidents">https://www.opendataphilly.org/dataset/crime-incidents</a>