

Experiment-10

Q CSV file manipulation

1 Import file crime_incidents.csv into Databricks File System.

2 Import SQLContext

→ from pyspark.sql import SQLContext

3 Load the file and create Dataframe.

→ crimeincidentDF = spark.read.csv("dbfs:/Filestore/tables/crime_incidents.csv", inferSchema="true", header="true")

4 Collect

→ crimeincidentDF.collect()

5 Take

→ crimeincidentDF.take(2)

6 Show

→ crimeincidentDF.show()

7 Shows first 2 records in tabular format.

→ crimeincidentDF.show(2)

8 Finding no. of records.

→ crimeincidentDF.count()

Topic
9. Printing the first row of a data set.
→ `crimeincidentDF.first()`

10 Schema of dataframe.
→ `primeincidentDF.printSchema()`

11. Describe dataframe.
→ crimeincidentDF.describe(), show()

12 Select function.

```
→ newDF = crimeincidentDF.select('IncidentNum?',  
                                  'category', 'descript')  
  
newDF.show()
```

13 Selecting all columns.
 → newDF = crimeincidentDF, select ("*", " ")
 newDF, show()

14 Correlation between 2 fields.
→ crimeincidentDF.corr('X', 'Y')

15 Remove Duplicate records.
→ newDF = crimeincidentDF.dropDuplicates()
newDF.show()

16 Finding name of a specific column using index.
→ `crimeincident [4]`

17 withColumn

```
→ crimeincidentNC = crimeincidentDF.withColumn(
  'IncidentNum', crimeincidentDF.IncidentNum).withColumn(
  'category', crimeincidentDF.category).withColumn(
  'Resolution', crimeincidentDF.Resolution)
crimeincidentNC.show()
```

18 Registering table as a template, and applying SQL queries to it.

```
→ crimeincidentNC.createOrReplaceTempView("crime incident
- table")
newDF = sqlContext.sql("select * from crime incident
table WHERE Resolution LIKE \'ARRE%-\'")
newDF.show()
```

19 Filtering Operations

```
→ subsetDF = crimeincidentNC.filter(crimeincidentNC.
Resolution == 'NONE')
subsetDF.show(10)
```

20 Filtering through multiple conditions

```
→ crimeincidentNC.filter(crimeincidentNC.Resolution
== 'LOCATED').filter(crimeincidentNC.category == 'ASSAULT')
crimeincidentNC.show()
```

21. Groupby

→ `crimeincidentDF.groupby(['category', crimeincidentDF.category]).count().show()`

22. Applying Aggregation functions after grouping.

→ `crimeincidentDF.groupby().avg().show()`

23. Aggregation using column 'Resolution'

→ `newDF = crimeincidentDF.groupby(crimeincidentDF.Resolution).newDF.count().show()`

24. Display 15 elements of dataframe.

→ `crimeincidentDF.show(15)`

25. Writing the dataframe as Parquet file.

→ `crimeincidentDF.write.parquet("dbfs:/Filestore/table/crime_incidents.parquet")`

26. The Parquet file will be stored by default in HDFS.

→ `newDF = spark.read.parquet("dbfs:/Filestore/table/crime_incidents.parquet", inferSchema = "true", header = "true")`
`newDF.show()`

Outputs:

Out[4]: Row(IncidentNum = 50436712, Category = 'ASSAULT',
 Descript = 'BATTERY', Dayofweek = 'Wednesday', Date = '04/20/2005',
 12:00 AM', Time = '04:00', Pddistrict = 'MISSION', Resolution = 'NONE',
 Address = '18th ST/CASTRO ST',)

Out[6]:

IncidentNum	Category	Descript	Dayofweek	Date	Time	Pddistrict
50436712	ASSAULT	BATTERY	Wednesday	04/20/05	04:00	MISSION
80049078	LARCENY/THEFT	GRAND THEFT	Sunday	01/13/8	18:00	PARK
		Resolution	Address	X	Y	Location
		NONE	18 th ST/CASTRO	-122.4350	32.760	(37.760...
		NONE	1100 Block of	-122.4467	32.762	(32.7625...

Out[8] crimeincidentDF.count()

: 1888567

Out[9]: Row(IncidentNum = 50436712, Category = 'ASSAULT',
 Description = 'BATTERY', Dayofweek = 'Wednesday', Date = '04/20/2005',
 12:00:00 AM', Time = '04:00', Pddistrict = 'MISSION', Resolution = 'NONE',
 Address = '18th ST/CASTRO ST', X = -122.4350, Y = 32.760882,
 Location = '(32.76091, -122.4350...)'

Out[10]: root.

- 1-- IncidentNum: integer (nullable = true)
- 1-- Category: string (nullable = true)
- 1-- Descript: string (nullable = true)
- 1-- DayOfWeek: string (nullable = true)
- 1-- Date: string (nullable = true)
- 1-- Time: string (nullable = true)
- 1-- PdDistrict: string (nullable = true)
- 1-- Resolution: string (nullable = true)
- 1-- Address: string (nullable = true)
- 1-- X: double (nullable = true)
- 1-- Y: double (nullable = true)
- 1-- Location: string (nullable = true)

Out[11]:

Summary	IncidentNum	Category	Descript	DayOfWeek	Date
count	1888567	1888567	1888567	1888567	1888567
mean	9.27955639	NULL	NULL	NULL	NULL
stddev	4.025831484	NULL	NULL	NULL	NULL
min	3949	ARSON	CHILD	Friday	01/01/06
max	991582322	WEAPON LAW	COURT	Wednesday	12/61/15

Out[12]:

IncidentNum	Category	Descript
50436212	ASSAULT	BATTERY
80049078	LARCENY/THEFT	GRAND THEFT AUTO
130366639	ASSAULT	AGGRAVATED ASSAULT

Out[14]: Correlation = 0.556556135469178

Out[16]: crimeincidentOF[4]
(column <'Date'>)

Out[18]:

IncidentNum	Category	Descript	Dayofweek	Date	Time
130366639	ASSAULT	AGGRAVATED	SUNDAY	05/05/13	04:10
30810835	DRIVING	DRIVING	TUESDAY	7/8/13	12:00

Pddistrict	Resolution	Address	X	Y	Location
INGLESIDE	ARREST, BOOKED	0 Block of STG	-122.442	32.72	(32.72...
SOUTHERN	ARREST, BOOKED	MAISON ST	-122.408	32.78	(32.78...

Out[19]:

IncidentNum	Category	Descript	Dayofweek	Date	Time	Pddistrict
50436712	ASSAULT	BATTERY	Wednesday	4/20/5	04:00	MISSION
80049078	ARCNEY	GRAND	Sunday	1/13/8	18:00	PARK

Resolution	Address	X	Y	Location
NONE	18th ST / CASTRO ST	-122.43500	32.760	(32.760,
NONE	1100 Block of CIA	-122.48500	32.762	(32.762,

Out[21]:

Category	Category	count
FRAUD	FRAUD	36004
SUICIDE	SUICIDE	1131
LIQVOR LAWS	LIQVOR LAWS	3860
SECONDARY CODES	SECONDARY CODES	21636

OUT[22]: Aggregation function

avg(IncidentNum)	avg(X)	avg(Y)
9.27955637558	122.42719103	32.741759506

OUT[23]:

Resolution	count
JUVENILE BOOKED	11789
EXCEPTIONAL CLEAR	3699
ARREST, BOOKED	445847
CLEARED-CONTACT	590

OUT[25]:

path dbfs:/Filestore/tables/crime_incidents.parquet already exists.