

Προχωρημένα Θέματα Βάσεων Δεδομένων

Μαρία Κοιλαλού | Μυρτώ Ορφανάκου
AM:03119211 | AM:

13 Ιανουαρίου 2024

1 Create DataFrame

```
1 schema1 = "`DR_NO` STRING, \  
2 `Date Rptd` STRING, \  
3 `DATE OCC` STRING, \  
4 `TIME OCC` INTEGER, \  
5 `AREA` INTEGER, \  
6 `AREA NAME` STRING, \  
7 `Rpt Dist No` INTEGER, \  
8 `Part 1-2` INTEGER, \  
9 `Crm Cd` INTEGER, \  
10 `Crm Cd Desc` STRING, \  
11 `Mocodes` STRING, \  
12 `Vict Age` INTEGER, \  
13 `Vict Sex` STRING, \  
14 `Vict Descent` STRING, \  
15 `Premis Cd` INTEGER, \  
16 `Premis Desc` STRING, \  
17 `Weapon Used Cd` INTEGER, \  
18 `Weapon Desc` STRING, \  
19 `Status` STRING, \  
20 `Status Desc` STRING, \  
21 `Crm Cd 1` INTEGER, \  
22 `Crm Cd 2` INTEGER, \  
23 `Crm Cd 3` INTEGER, \  
24 `Crm Cd 4` INTEGER, \  
25 `LOCATION` STRING, \  
26 `Cross Street` STRING, \  
27 `LAT` DOUBLE, \  
28 `LON` DOUBLE"  
29  
30  
31 data1 = spark.read.csv("/user/ubuntu/ta/advanced-db/data/crime_data_2010.csv", header=True, schema=schema1)  
32 data2 = spark.read.csv("/user/ubuntu/ta/advanced-db/data/crime_data_2020.csv", header=True, schema=schema1)  
33 df = data1.union(data2).distinct()  
34  
35 df = df.withColumn("Date Rptd", to_date(col("Date Rptd"), "MM/dd/yyyy hh:mm:ss a")) \  
36 .withColumn("DATE OCC", to_date(col("DATE OCC"), "MM/dd/yyyy hh:mm:ss a"))  
37  
38 df.count()  
39 print(f"Total number of rows: {df.count()}")  
40  
41 df.printSchema()
```

```
Total number of rows: 2913595  
root  
|-- DR_NO: string (nullable = true)
```

```
|-- Date Rptd: date (nullable = true)
|-- DATE OCC: date (nullable = true)
|-- TIME OCC: integer (nullable = true)
|-- AREA: integer (nullable = true)
|-- AREA NAME: string (nullable = true)
|-- Rpt Dist No: integer (nullable = true)
|-- Part 1-2: integer (nullable = true)
|-- Crm Cd: integer (nullable = true)
|-- Crm Cd Desc: string (nullable = true)
|-- Mocodes: string (nullable = true)
|-- Vict Age: integer (nullable = true)
|-- Vict Sex: string (nullable = true)
|-- Vict Descent: string (nullable = true)
|-- Premis Cd: integer (nullable = true)
|-- Premis Desc: string (nullable = true)
|-- Weapon Used Cd: integer (nullable = true)
|-- Weapon Desc: string (nullable = true)
|-- Status: string (nullable = true)
|-- Status Desc: string (nullable = true)
|-- Crm Cd 1: integer (nullable = true)
|-- Crm Cd 2: integer (nullable = true)
|-- Crm Cd 3: integer (nullable = true)
|-- Crm Cd 4: integer (nullable = true)
|-- LOCATION: string (nullable = true)
|-- Cross Street: string (nullable = true)
|-- LAT: double (nullable = true)
|-- LON: double (nullable = true)
```