

DB0201EN-Week4-2-2-PeerAssign-v5-py

April 22, 2019

Assignment: Notebook for Peer Assignment

1 Introduction

Using this Python notebook you will: 1. Understand 3 Chicago datasets

1. Load the 3 datasets into 3 tables in a Db2 database 1. Execute SQL queries to answer assignment questions

1.1 Understand the datasets

To complete the assignment problems in this notebook you will be using three datasets that are available on the city of Chicago's Data Portal: 1. Socioeconomic Indicators in Chicago 1. Chicago Public Schools 1. Chicago Crime Data

1.1.1 1. Socioeconomic Indicators in Chicago

This dataset contains a selection of six socioeconomic indicators of public health significance and a "hardship index," for each Chicago community area, for the years 2008 – 2012.

For this assignment you will use a snapshot of this dataset which can be downloaded from: <https://ibm.box.com/shared/static/05c3415cbfbtfnr2fx4atenb2sd361ze.csv>

A detailed description of this dataset and the original dataset can be obtained from the Chicago Data Portal at: <https://data.cityofchicago.org/Health-Human-Services/Census-Data-Selected-socioeconomic-indicators-in-C/kn9c-c2s2>

1.1.2 2. Chicago Public Schools

This dataset shows all school level performance data used to create CPS School Report Cards for the 2011-2012 school year. This dataset is provided by the city of Chicago's Data Portal.

For this assignment you will use a snapshot of this dataset which can be downloaded from: <https://ibm.box.com/shared/static/f9gijv1gjmxxzycdhplzt01qtz0s7ew7.csv>

A detailed description of this dataset and the original dataset can be obtained from the Chicago Data Portal at: <https://data.cityofchicago.org/Education/Chicago-Public-Schools-Progress-Report-Cards-2011-/9xs2-f89t>

1.1.3 3. Chicago Crime Data

This dataset reflects reported incidents of crime (with the exception of murders where data exists for each victim) that occurred in the City of Chicago from 2001 to present, minus the most recent seven days.

This dataset is quite large - over 1.5GB in size with over 6.5 million rows. For the purposes of this assignment we will use a much smaller sample of this dataset which can be downloaded from: <https://ibm.box.com/shared/static/svflyugsr9zbqy5bmowgswqemfpm1x7f.csv>

A detailed description of this dataset and the original dataset can be obtained from the Chicago Data Portal at: <https://data.cityofchicago.org/Public-Safety/Crimes-2001-to-present/ijzp-q8t2>

1.1.4 Download the datasets

In many cases the dataset to be analyzed is available as a .CSV (comma separated values) file, perhaps on the internet. Click on the links below to download and save the datasets (.CSV files):

1. **CENSUS_DATA**: <https://ibm.box.com/shared/static/05c3415cbfbtfnr2fx4atenb2sd361ze.csv>

1. **CHICAGO_PUBLIC_SCHOOLS** <https://ibm.box.com/shared/static/f9gjvj1gjmxxzycdhplzt01qtz0s7ew7.csv>

1. **CHICAGO_CRIME_DATA**: <https://ibm.box.com/shared/static/svflyugsr9zbqy5bmowgswqemfpm1x7f.csv>

NOTE: Ensure you have downloaded the datasets using the links above instead of directly from the Chicago Data Portal. The versions linked here are subsets of the original datasets and have some of the column names modified to be more database friendly which will make it easier to complete this assignment.

1.1.5 Store the datasets in database tables

To analyze the data using SQL, it first needs to be stored in the database.

While it is easier to read the dataset into a Pandas dataframe and then PERSIST it into the database as we saw in Week 3 Lab 3, it results in mapping to default datatypes which may not be optimal for SQL querying. For example a long textual field may map to a CLOB instead of a VARCHAR.

Therefore, **it is highly recommended to manually load the table using the database console LOAD tool, as indicated in Week 2 Lab 1 Part II.** The only difference with that lab is that in Step 5 of the instructions you will need to click on create "(+) New Table" and specify the name of the table you want to create and then click "Next".

Now open the Db2 console, open the LOAD tool, Select / Drag the .CSV file for the first dataset, Next create a New Table, and then follow the steps on-screen instructions to load the data. Name the new tables as follows:

1. **CENSUS_DATA**
2. **CHICAGO_PUBLIC_SCHOOLS**
3. **CHICAGO_CRIME_DATA**

1.1.6 Connect to the database

Let us first load the SQL extension and establish a connection with the database

```
In [1]: %load_ext sql
```

In the next cell enter your db2 connection string. Recall you created Service Credentials for your Db2 instance in first lab in Week 3. From the **uri** field of your Db2 service credentials copy everything after db2:// (except the double quote at the end) and paste it in the cell below after ibm_db_sa://

```
In [2]: # Remember the connection string is of the format:
        # %sql ibm_db_sa://my-username:my-password@my-hostname:my-port/my-db-name
        # Enter the connection string for your Db2 on Cloud database instance below
        %sql ibm_db_sa://lln32654:69d16lp9%5Ecw5lj4g@dashdb-txn-sbox-yp-dal09-03.services.dal.ibm.com:50000/BLUDB
```

```
Out[2]: 'Connected: lln32654@BLUDB'
```

1.2 Problems

Now write and execute SQL queries to solve assignment problems

1.2.1 Problem 1

Find the total number of crimes recorded in the CRIME table

```
In [40]: # Rows in Crime table
         %sql select count(*) AS "Number_of_Crimes" from chicago_crime_data
```

```
* ibm_db_sa://lln32654:***@dashdb-txn-sbox-yp-dal09-03.services.dal.ibm.com:50000/BLUDB
Done.
```

```
Out[40]: [(Decimal('533'),)]
```

1.2.2 Problem 2

Retrieve first 10 rows from the CRIME table

```
In [4]: %sql select * from chicago_crime_data \
        fetch first 10 rows only
```

```
* ibm_db_sa://lln32654:***@dashdb-txn-sbox-yp-dal09-03.services.dal.ibm.com:50000/BLUDB
Done.
```

```
Out[4]: [(3512276, 'HK587712', '08/28/2004 05:50:56 PM', '047XX S KEDZIE AVE', '890', 'THEFT', 'FROM BU
(3406613, 'HK456306', '06/26/2004 12:40:00 PM', '009XX N CENTRAL PARK AVE', '820', 'THEFT', '$5
(8002131, 'HT233595', '04/04/2011 05:45:00 AM', '043XX S WABASH AVE', '820', 'THEFT', '$500 AND
(7903289, 'HT133522', '12/30/2010 04:30:00 PM', '083XX S KINGSTON AVE', '840', 'THEFT', 'FINANC
(10402076, 'HZ138551', '02/02/2016 07:30:00 PM', '033XX W 66TH ST', '820', 'THEFT', '$500 AND UN
(7732712, 'HS540106', '09/29/2010 07:59:00 AM', '006XX W CHICAGO AVE', '810', 'THEFT', 'OVER $5
(10769475, 'HZ534771', '11/30/2016 01:15:00 AM', '050XX N KEDZIE AVE', '810', 'THEFT', 'OVER $50
(4494340, 'HL793243', '12/16/2005 04:45:00 PM', '005XX E PERSHING RD', '860', 'THEFT', 'RETAIL T
(3778925, 'HL149610', '01/28/2005 05:00:00 PM', '100XX S WASHTENAW AVE', '810', 'THEFT', 'OVER
(3324217, 'HK361551', '05/13/2004 02:15:00 PM', '033XX W BELMONT AVE', '820', 'THEFT', '$500 AN
```

1.2.3 Problem 3

How many crimes involve an arrest?

```
In [41]: %sql select count(*) AS "Number_of_Crimes_Involving_an_Arrest" from chicago_crime_data where a
* ibm_db_sa://lln32654:***@dashdb-txn-sbox-yp-dal09-03.services.dal.ibm.com:50000/BLUDB
Done.
```

```
Out[41]: [(Decimal('163'),)]
```

1.2.4 Problem 4

Which unique types of crimes have been recorded at GAS STATION locations?

```
In [6]: %sql select distinct primary_type, location_description from chicago_crime_data where location_descrip
* ibm_db_sa://lln32654:***@dashdb-txn-sbox-yp-dal09-03.services.dal.ibm.com:50000/BLUDB
Done.
```

```
Out[6]: [('CRIMINAL TRESPASS', 'GAS STATION'),
         ('NARCOTICS', 'GAS STATION'),
         ('ROBBERY', 'GAS STATION'),
         ('THEFT', 'GAS STATION')]
```

Hint: Which column lists types of crimes e.g. THEFT?

1.2.5 Problem 5

In the CENUS_DATA table list all Community Areas whose names start with the letter 'B'.

```
In [7]: %sql select community_area_name from census_data where census_data.community_area_name LIKE 'B%'
* ibm_db_sa://lln32654:***@dashdb-txn-sbox-yp-dal09-03.services.dal.ibm.com:50000/BLUDB
Done.
```

```
Out[7]: [('Belmont Cragin'),
         ('Burnside'),
         ('Brighton Park'),
         ('Bridgeport'),
         ('Beverly')]
```

1.2.6 Problem 6

Which schools in Community Areas 10 to 15 are healthy school certified?

```
In [33]: %sql select name_of_school, community_area_number, healthy_school_certified from CHICAGO_PUBLIC_SCHOOLS
* ibm_db_sa://lln32654:***@dashdb-txn-sbox-yp-dal09-03.services.dal.ibm.com:50000/BLUDB
Done.
```

```
Out[33]: [('Rufus M Hitch Elementary School', 10, 'Yes')]
```

1.2.7 Problem 7

What is the average school Safety Score?

```
In [50]: %sql select round(avg(safety_score), 2) AS "Average_Safety_Score" from CHICAGO_PUBLIC_SCHOOLS
* ibm_db_sa://ln32654:***@dashdb-txn-sbox-yp-dal09-03.services.dal.ibm.com:50000/BLUDB
Done.
```

```
Out[50]: [(Decimal('49.500000'),)]
```

1.2.8 Problem 8

List the top 5 Community Areas by average College Enrollment [number of students]

```
In [87]: %sql select COMMUNITY_AREA_NAME, round(avg(COLLEGE_ENROLLMENT),2) as Average_College_Enrollment
group by COMMUNITY_AREA_NAME \
ORDER BY Average_College_Enrollment DESC\
fetch first 5 rows only
* ibm_db_sa://ln32654:***@dashdb-txn-sbox-yp-dal09-03.services.dal.ibm.com:50000/BLUDB
Done.
```

```
Out[87]: [('ARCHER HEIGHTS', Decimal('2411.500000')),
('MONTCLARE', Decimal('1317.000000')),
('WEST ELSDON', Decimal('1233.330000')),
('BRIGHTON PARK', Decimal('1205.880000')),
('BELMONT CRAGIN', Decimal('1198.830000'))]
```

1.2.9 Problem 9

Use a sub-query to determine which Community Area has the least value for school Safety Score?

```
In [92]: %sql select COMMUNITY_AREA_NAME, SAFETY_SCORE from CHICAGO_PUBLIC_SCHOOLS
SAFETY_SCORE = (SELECT MIN(SAFETY_SCORE) FROM CHICAGO_PUBLIC_SCHOOLS)
* ibm_db_sa://ln32654:***@dashdb-txn-sbox-yp-dal09-03.services.dal.ibm.com:50000/BLUDB
Done.
```

```
Out[92]: [('WASHINGTON PARK', 1)]
```

1.2.10 Problem 10

[Without using an explicit JOIN operator] Find the Per Capita Income of the Community Area which has a school Safety Score of 1.

```
In [196]: #
          %sql select community_area_number,PER_CAPITA_INCOME\
          from CENSUS_DATA where community_area_number in\
          (select community_area_number from CHICAGO_PUBLIC_SCHOOLS where SAFETY_SCORE=1)

* ibm_db_sa://lln32654:***@dashdb-txn-sbox-yp-dal09-03.services.dal.ibmcloud.net:50000/BLUDB
Done.
```

```
Out[196]: [(40, 13785)]
```

Copyright © 2018 cognitiveclass.ai. This notebook and its source code are released under the terms of the [MIT License](https://creativecommons.org/licenses/by/4.0/).