



Analyzing a real world data-set with SQL and Python

Estimated time needed: **15** minutes

Objectives

After completing this lab you will be able to:

- Understand a dataset of selected socioeconomic indicators in Chicago
- Learn how to store data in an SQLite database.
- Solve example problems to practice your SQL skills

Selected Socioeconomic Indicators in Chicago

The city of Chicago released a dataset of socioeconomic data to the Chicago City Portal. 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.

Scores on the hardship index can range from 1 to 100, with a higher index number representing a greater level of hardship.

A detailed description of the dataset can be found on [the city of Chicago's website](#), but to summarize, the dataset has the following variables:

- **Community Area Number** (`ca`): Used to uniquely identify each row of the dataset
- **Community Area Name** (`community_area_name`): The name of the region in the city of Chicago
- **Percent of Housing Crowded** (`percent_of_housing_crowded`): Percent of occupied housing units with more than one person per room
- **Percent Households Below Poverty** (`percent_households_below_poverty`): Percent of households living below the federal poverty line
- **Percent Aged 16+ Unemployed** (`percent_aged_16_unemployed`): Percent of persons over the age of 16 years that are unemployed

- **Percent Aged 25+ without High School Diploma** (`percent_aged_25_without_high_school_diploma`): Percent of persons over the age of 25 years without a high school education
- **Percent Aged Under 18 or Over 64**: Percent of population under 18 or over 64 years of age (`percent_aged_under_18_or_over_64`): (ie. dependents)
- **Per Capita Income** (`per_capita_income_`): Community Area per capita income is estimated as the sum of tract-level aggregate incomes divided by the total population
- **Hardship Index** (`hardship_index`): Score that incorporates each of the six selected socioeconomic indicators

In this Lab, we'll take a look at the variables in the socioeconomic indicators dataset and do some basic analysis with Python.

Connect to the database

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

The syntax for connecting to magic sql using sqlite is

%sql sqlite://DatabaseName

where DatabaseName will be your **.db** file

```
In [1]: !pip install ipython-sql  
%load_ext sql
```

```

Collecting ipython-sql
  Downloading ipython_sql-0.5.0-py3-none-any.whl.metadata (17 kB)
Collecting prettytable (from ipython-sql)
  Downloading prettytable-3.11.0-py3-none-any.whl.metadata (30 kB)
Requirement already satisfied: ipython in /opt/conda/lib/python3.11/site-packages (from ipython-sql) (8.22.2)
Requirement already satisfied: sqlalchemy>=2.0 in /opt/conda/lib/python3.11/site-packages (from ipython-sql) (2.0.30)
Collecting sqlparse (from ipython-sql)
  Downloading sqlparse-0.5.1-py3-none-any.whl.metadata (3.9 kB)
Requirement already satisfied: six in /opt/conda/lib/python3.11/site-packages (from ipython-sql) (1.16.0)
Requirement already satisfied: ipython-genutils in /opt/conda/lib/python3.11/site-packages (from ipython-sql) (0.2.0)
Requirement already satisfied: typing-extensions>=4.6.0 in /opt/conda/lib/python3.11/site-packages (from sqlalchemy>=2.0->ipython-sql) (4.11.0)
Requirement already satisfied: greenlet!=0.4.17 in /opt/conda/lib/python3.11/site-packages (from sqlalchemy>=2.0->ipython-sql) (3.0.3)
Requirement already satisfied: decorator in /opt/conda/lib/python3.11/site-packages (from ipython->ipython-sql) (5.1.1)
Requirement already satisfied: jedi>=0.16 in /opt/conda/lib/python3.11/site-packages (from ipython->ipython-sql) (0.19.1)
Requirement already satisfied: matplotlib-inline in /opt/conda/lib/python3.11/site-packages (from ipython->ipython-sql) (0.1.7)
Requirement already satisfied: prompt-toolkit<3.1.0,>=3.0.41 in /opt/conda/lib/python3.11/site-packages (from ipython->ipython-sql) (3.0.42)
Requirement already satisfied: pygments>=2.4.0 in /opt/conda/lib/python3.11/site-packages (from ipython->ipython-sql) (2.18.0)
Requirement already satisfied: stack-data in /opt/conda/lib/python3.11/site-packages (from ipython->ipython-sql) (0.6.2)
Requirement already satisfied: traitlets>=5.13.0 in /opt/conda/lib/python3.11/site-packages (from ipython->ipython-sql) (5.14.3)
Requirement already satisfied: pexpect>4.3 in /opt/conda/lib/python3.11/site-packages (from ipython->ipython-sql) (4.9.0)
Requirement already satisfied: wcwidth in /opt/conda/lib/python3.11/site-packages (from prettytable->ipython-sql) (0.2.13)
Requirement already satisfied: parso<0.9.0,>=0.8.3 in /opt/conda/lib/python3.11/site-packages (from jedi>=0.16->ipython->ipython-sql) (0.8.4)
Requirement already satisfied: ptyprocess>=0.5 in /opt/conda/lib/python3.11/site-packages (from pexpect>4.3->ipython->ipython-sql) (0.7.0)
Requirement already satisfied: executing>=1.2.0 in /opt/conda/lib/python3.11/site-packages (from stack-data->ipython->ipython-sql) (2.0.1)
Requirement already satisfied: asttokens>=2.1.0 in /opt/conda/lib/python3.11/site-packages (from stack-data->ipython->ipython-sql) (2.4.1)
Requirement already satisfied: pure-eval in /opt/conda/lib/python3.11/site-packages (from stack-data->ipython->ipython-sql) (0.2.2)
Downloading ipython_sql-0.5.0-py3-none-any.whl (20 kB)
Downloading prettytable-3.11.0-py3-none-any.whl (28 kB)
Downloading sqlparse-0.5.1-py3-none-any.whl (44 kB)

```

44.2/44.2 kB 7.4 MB/s eta 0:00:00

Installing collected packages: sqlparse, prettytable, ipython-sql

Successfully installed ipython-sql-0.5.0 prettytable-3.11.0 sqlparse-0.5.1

In [2]: `import csv, sqlite3`

```
con = sqlite3.connect("socioeconomic.db")
```

```

cur = con.cursor()
!pip install pandas

Collecting pandas
  Downloading pandas-2.2.3-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (89 kB)
   ━━━━━━━━━━━━━━━━ 89.9/89.9 kB 10.9 MB/s eta 0:00:00

Collecting numpy>=1.23.2 (from pandas)
  Downloading numpy-2.1.2-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl.metadata (60 kB)
   ━━━━━━━━━━━━━━ 60.9/60.9 kB 8.6 MB/s eta 0:00:00

Requirement already satisfied: python-dateutil>=2.8.2 in /opt/conda/lib/python3.11/site-packages (from pandas) (2.9.0)
Requirement already satisfied: pytz>=2020.1 in /opt/conda/lib/python3.11/site-packages (from pandas) (2024.1)
Collecting tzdata>=2022.7 (from pandas)
  Downloading tzdata-2024.2-py2.py3-none-any.whl.metadata (1.4 kB)
Requirement already satisfied: six>=1.5 in /opt/conda/lib/python3.11/site-packages (from python-dateutil>=2.8.2->pandas) (1.16.0)
  Downloading pandas-2.2.3-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (13.1 MB)
   ━━━━━━━━━━━━━━ 13.1/13.1 kB 90.8 MB/s eta 0:00:00:00:01
0:01
  Downloading numpy-2.1.2-cp311-cp311-manylinux_2_17_x86_64.manylinux2014_x86_64.whl (16.3 kB)
   ━━━━━━━━━━━━━━ 16.3/16.3 kB 85.7 MB/s eta 0:00:00:00:01
0:01
  Downloading tzdata-2024.2-py2.py3-none-any.whl (346 kB)
   ━━━━━━━━━━━━ 346.6/346.6 kB 43.4 MB/s eta 0:00:00

Installing collected packages: tzdata, numpy, pandas
Successfully installed numpy-2.1.2 pandas-2.2.3 tzdata-2024.2

```

In [3]: `%sql sqlite:///socioeconomic.db`

Store the dataset in a Table

In many cases the dataset to be analyzed is available as a .CSV (comma separated values) file, perhaps on the internet. To analyze the data using SQL, it first needs to be stored in the database.

We will first read the csv files from the given url into pandas dataframes

Next we will be using the `df.to_sql()` function to convert each csv file to a table in sqlite with the csv data loaded in it.

In [4]: `import pandas
df = pandas.read_csv('https://data.cityofchicago.org/resource/jcxq-k9xf.csv')
df.to_sql("chicago_socioeconomic_data", con, if_exists='replace', index=False, method`

Out[4]: 78

You can verify that the table creation was successful by making a basic query like:

In [5]: `%sql SELECT * FROM chicago_socioeconomic_data limit 5;`

```
* sqlite:///socioeconomic.db
Done.
```

Out[5]:

ca	community_area_name	percent_of_housing_crowded	percent_households_below_poverty
1.0	Rogers Park	7.7	23.6
2.0	West Ridge	7.8	17.2
3.0	Uptown	3.8	24.0
4.0	Lincoln Square	3.4	10.9
5.0	North Center	0.3	7.5

Problems

Problem 1

How many rows are in the dataset?

In [6]:

```
%sql SELECT COUNT(*) FROM chicago_socioeconomic_data;
```

```
* sqlite:///socioeconomic.db
Done.
```

Out[6]:

COUNT(*)
78

► Click here for the solution

Problem 2

How many community areas in Chicago have a hardship index greater than 50.0?

In [8]:

```
%sql select count(ca) from chicago_socioeconomic_data where hardship_index>50.0;
```

```
* sqlite:///socioeconomic.db
Done.
```

Out[8]:

count(ca)
38

► Click here for the solution

Problem 3

What is the maximum value of hardship index in this dataset?

In [9]:

```
%sql select max(hardship_index) from chicago_socioeconomic_data;
```

```
* sqlite:///socioeconomic.db
Done.
```

Out[9]: **max(hardship_index)**

98.0

► Click here for the solution

Problem 4

Which community area which has the highest hardship index?

In [10]: **%sql select community_area_name from chicago_socioeconomic_data where hardship_inde**

```
* sqlite:///socioeconomic.db
Done.
```

Out[10]: **community_area_name**

Riverdale

► Click here for the solution

Problem 5

Which Chicago community areas have per-capita incomes greater than \$60,000?

In [11]: **%sql SELECT community_area_name FROM chicago_socioeconomic_data WHERE per_capita_in**

```
* sqlite:///socioeconomic.db
Done.
```

Out[11]: **community_area_name**

Lake View

Lincoln Park

Near North Side

Loop

► Click here for the solution

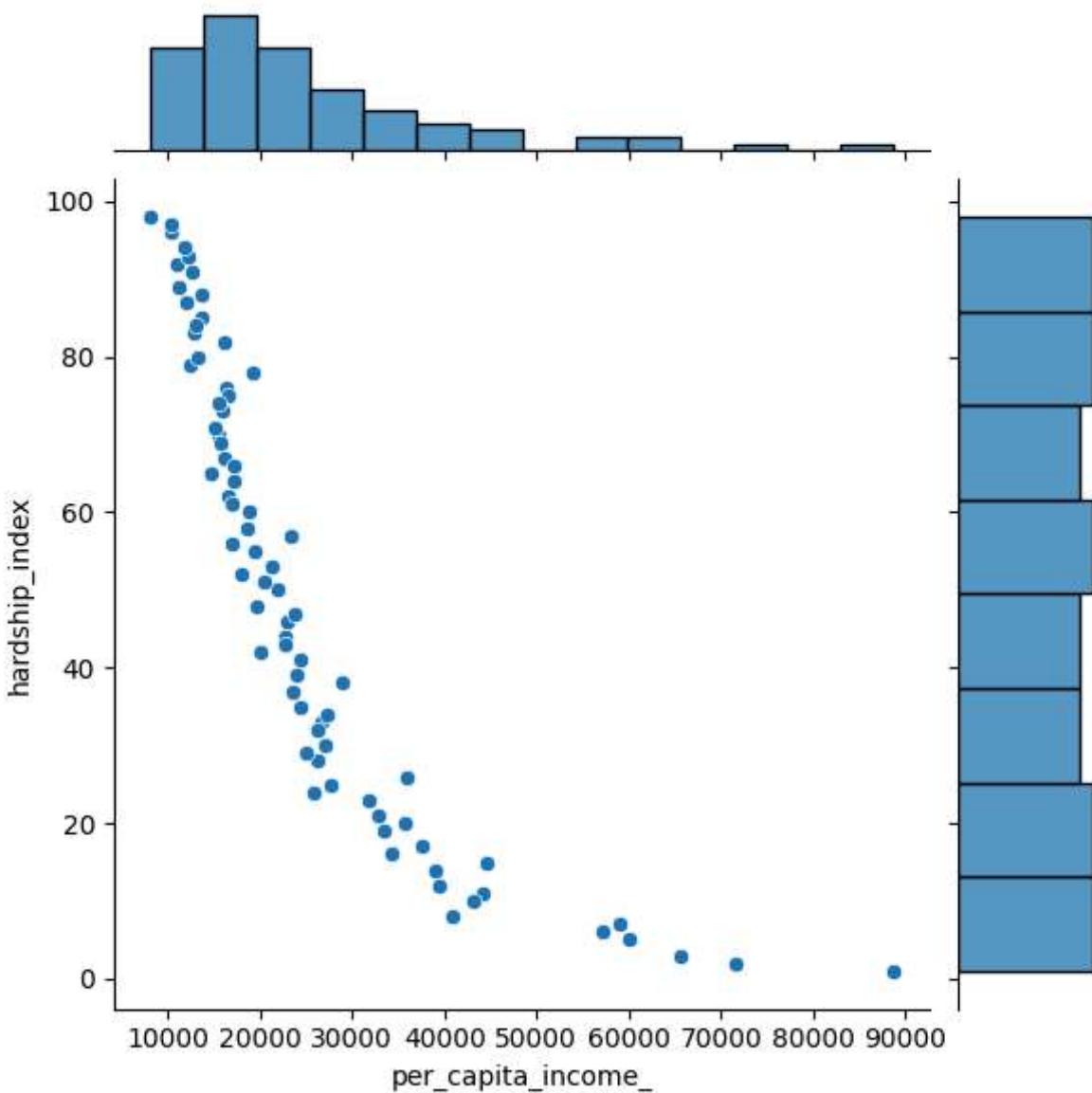
Problem 6

Create a scatter plot using the variables `per_capita_income_` and `hardship_index`. Explain the correlation between the two variables.

In [16]: **!pip install matplotlib
#import matplotlib.pyplot as plt
%matplotlib inline
!pip install seaborn
import seaborn as sns**

```
income_vs_hardship = %sql SELECT per_capita_income_, hardship_index FROM chicago_so
plot = sns.jointplot(x='per_capita_income_',y='hardship_index', data=income_vs_hard
```

```
Requirement already satisfied: matplotlib in /opt/conda/lib/python3.11/site-packages  
(3.9.2)  
Requirement already satisfied: contourpy>=1.0.1 in /opt/conda/lib/python3.11/site-packages (from matplotlib) (1.3.0)  
Requirement already satisfied: cycler>=0.10 in /opt/conda/lib/python3.11/site-packages (from matplotlib) (0.12.1)  
Requirement already satisfied: fonttools>=4.22.0 in /opt/conda/lib/python3.11/site-packages (from matplotlib) (4.54.1)  
Requirement already satisfied: kiwisolver>=1.3.1 in /opt/conda/lib/python3.11/site-packages (from matplotlib) (1.4.7)  
Requirement already satisfied: numpy>=1.23 in /opt/conda/lib/python3.11/site-packages (from matplotlib) (2.1.2)  
Requirement already satisfied: packaging>=20.0 in /opt/conda/lib/python3.11/site-packages (from matplotlib) (24.0)  
Requirement already satisfied: pillow>=8 in /opt/conda/lib/python3.11/site-packages (from matplotlib) (11.0.0)  
Requirement already satisfied: pyparsing>=2.3.1 in /opt/conda/lib/python3.11/site-packages (from matplotlib) (3.2.0)  
Requirement already satisfied: python-dateutil>=2.7 in /opt/conda/lib/python3.11/site-packages (from matplotlib) (2.9.0)  
Requirement already satisfied: six>=1.5 in /opt/conda/lib/python3.11/site-packages (from python-dateutil>=2.7->matplotlib) (1.16.0)  
Requirement already satisfied: seaborn in /opt/conda/lib/python3.11/site-packages (0.13.2)  
Requirement already satisfied: numpy!=1.24.0,>=1.20 in /opt/conda/lib/python3.11/site-packages (from seaborn) (2.1.2)  
Requirement already satisfied: pandas>=1.2 in /opt/conda/lib/python3.11/site-packages (from seaborn) (2.2.3)  
Requirement already satisfied: matplotlib!=3.6.1,>=3.4 in /opt/conda/lib/python3.11/site-packages (from seaborn) (3.9.2)  
Requirement already satisfied: contourpy>=1.0.1 in /opt/conda/lib/python3.11/site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (1.3.0)  
Requirement already satisfied: cycler>=0.10 in /opt/conda/lib/python3.11/site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (0.12.1)  
Requirement already satisfied: fonttools>=4.22.0 in /opt/conda/lib/python3.11/site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (4.54.1)  
Requirement already satisfied: kiwisolver>=1.3.1 in /opt/conda/lib/python3.11/site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (1.4.7)  
Requirement already satisfied: packaging>=20.0 in /opt/conda/lib/python3.11/site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (24.0)  
Requirement already satisfied: pillow>=8 in /opt/conda/lib/python3.11/site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (11.0.0)  
Requirement already satisfied: pyparsing>=2.3.1 in /opt/conda/lib/python3.11/site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (3.2.0)  
Requirement already satisfied: python-dateutil>=2.7 in /opt/conda/lib/python3.11/site-packages (from matplotlib!=3.6.1,>=3.4->seaborn) (2.9.0)  
Requirement already satisfied: pytz>=2020.1 in /opt/conda/lib/python3.11/site-packages (from pandas>=1.2->seaborn) (2024.1)  
Requirement already satisfied: tzdata>=2022.7 in /opt/conda/lib/python3.11/site-packages (from pandas>=1.2->seaborn) (2024.2)  
Requirement already satisfied: six>=1.5 in /opt/conda/lib/python3.11/site-packages (from python-dateutil>=2.7->matplotlib!=3.6.1,>=3.4->seaborn) (1.16.0)  
* sqlite:///socioeconomic.db  
Done.
```



► Click here for the solution

Conclusion

Now that you know how to do basic exploratory data analysis using SQL and python visualization tools, you can further explore this dataset to see how the variable `per_capita_income_` is related to `percent_households_below_poverty` and `percent_aged_16_unemployed`. Try to create interesting visualizations!

Summary

In this lab you learned how to store a real world data set from the internet in a database, gain insights into data using SQL queries. You also visualized a portion of the data in the database to see what story it tells.

Author

Rav Ahuja

{toggle}

{toggle}|

{toggle}

{toggle}

{toggle}|

{toggle}|

{toggle}|

© IBM Corporation 2020. All rights reserved.