

Week 4 - Declarative SQL Analysis

Where do we get data from??

- provided (csv, excel, xml, json) datasets.

by organisation / available from online

- web scraping

Relational databases

relation: named, two-dimensional table data

This week's objective

RDBMS + Python X analysis with SQL

CSV → DB via python

then use SQL to analyse that data.

SQL - structured query language.

• DDL, DML

create, drop
update

retrieve information

alter
regarding tables

Insert, update, ~~add~~ delete, select from

Table constraints & relational key

- ~~pk~~ unique, not null

→ also check(), default

- FK: refers to candidate key of parent table.

} can be
composite

{ DML stuffs

Loading data into SQL

Approach A: \copy

(in psql)

\COPY {table} FROM {filename} CSV [HEADER]

pro: fast, no programming

con: only supports 1:1 mapping (from only one file source)

no data cleaning / transformation \therefore prone to error, crash

* Approach B: Python

1. create connection (pgconnect)

2. create tables (conn.execute)

3. df = pd.read_csv()

4. df.to_sql on conn

manually-defined function

pro: flexible data cleaning / transformation

con: coding required

df.to_sql(table-name, conn, if-exists = 'replace')

Data cleaning with Python

• type conversion

- int(), float()

- datetime.strptime()

• missing values

- to None or NaN

For every row in a specified table,

- check row is one of ~~top~~ special values

- check casting runs w/o error

→ if error, default value.

skip np.nan