

sqawk

or

CLI SQL for CSV:-)

i.e.

making some shell tasks on CSV files easier

Shell one-liners

a. `$ awk '$3 > 12.5' < myfile`

b. `$ join file1 file2`

supposing both files are sorted on the join field, in this case #1

Similar tasks

- a'. extract all rows of myfile where the value in column 3 is above average
- b'. join file1 with file2, but on a composite join field (e.g. hospital ID + patient ID)
- c'. join more than two files (e.g. genotypes, covariates, eigenvalues)

As Database Operations

... you can express them as a single SQL query:

a'. `SELECT * FROM t1 WHERE f3 >
(SELECT avg(f3));`

b'. `SELECT t1.* FROM t1 JOIN t2 USING
(f1,f2,f3);`

c'. `SELECT ... FROM t1 JOIN t2 ...
JOIN t3 ...;`

where tables (t1, etc.) contain files, columns (f1, etc.) contain file fields

but...

to do SQL queries on a file, you need:

- to create a database
- to create a table for the file
- to import the file's data into the table
- (usually) a SQL server

...which is a bit unwieldy.

Wish List

the successful candidate will:

- automatically create a database
- automatically create db tables from CSV files
- automatically import content into the tables
- run a SQL query
- print out the result
- be a shell filter

SQLite

- C library (not server)
- small, fast

→ if we can automate table creation and population, we're done.

Syntax

```
$ sqawk [opts] ([file opts] file)+  
SQL
```

Selected options:

- **-i**: specifies index fields
- **-p**: specifies primary key
- **-q**: shows the generated SQL

Examples

a'. `$ sqawk myfile.csv 'SELECT *
FROM myfile WHERE f3 > (SELECT
avg(f3))'`

b'. `sqawk file1.csv file2.csv 'SELECT
* FROM file1 JOIN file2 USING
(f1,f2,f3);`

Checking for valid IDs

- file valid: list of valid IDs
- file dubious: uncertain IDs (among other data)

```
$ ./sqawk valid dubious 'SELECT *  
FROM dubious WHERE dubious.spc NOT IN  
(SELECT spc FROM valid)'
```


In an ideal world:

- data is consistently formatted
- data formats are compatible
- data is validated before use
- ...

In the real world...

- data is messy
- there is a plethora of incompatible formats
- nobody checks the data before sending it :-)
- ...

what can be done?

- export to CSV
- write code to systematically check the data
- ...

⇒ this is where sqawk might help.

