

Unlearning SQL

Leverage the SQL design patterns in Python

S.Lott

<https://fosstodon.org/@slott56>

<https://github.com/slott56>

14-May-2025

Table of Contents

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

- 1 SQL Overuse
- 2 SQL Design Patterns
- 3 Python Implementation Patterns
- 4 About Those Join Algorithms
- 5 Group By and Having — The Good Stuff
- 6 Conclusion

SQL is Helpful

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

- Folks know SQL.
- They are fluent in SQL's design patterns.
- They find it hard to convert SQL designs to Python.

This talk should help clarify SQL from a Python perspective

Database Overuse

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

Given a problem requiring a SQL-like summary.

- 1 Define tables
- 2 Quick load script
- 3 SQL SELECT

Easy, right?

Database Overuse

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

Given a problem requiring a SQL-like summary.

- 1 Define tables
- 2 Quick load script
- 3 SQL SELECT

Easy, right?

Maybe not

SQL Overheads

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

The database engine has overheads

Lots of them.

Locking. Storage management. Permissions. Serialization.

War Story:

Developer struggling with transient data processing.

The app had repeated **Create-Load-Query-Drop** cycles.

The Drop (it turns out) is both unpredictable and slow.

(Even if it's SQLite, there are overheads.)

How do we unlearn SQL?

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

Two steps to moving past SQL:

- ① Understand the SQL design patterns.
- ② Rework those design elements in Python.

SQL Design Patterns

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

Consider the core Select statement:

SELECT *expr*, ...

FROM *table*, ...

WHERE *condition*

We'll get to GROUP BY and HAVING later.

SELECT works like this

Unlearning
SQL

S.Lott

SQL Overuse

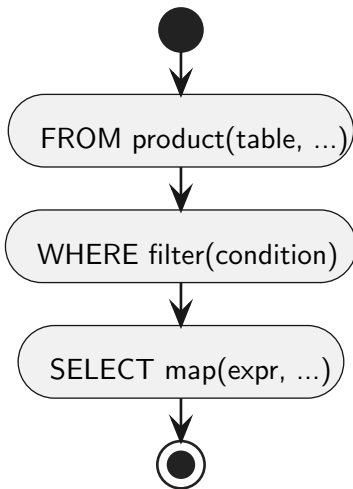
SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion



SELECT in Python

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

```
FROM t1, t2, ...
```

```
from_ = itertools.product(t1, t2, ...)
```

```
WHERE c
```

```
where = (row_tuple  
         for row_tuple in from_  
         if c(row_tuple))
```

```
SELECT ex1, ex2, ...
```

```
result = list(  
    (ex1(row_tuple), ex2(row_tuple), ...)   
    for row_tuple in where)
```

Good and Bad

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

- The Python code matches the SQL.
- A lot more syntax: `WHERE expr becomes (r for r in from_ if expr),`
- `SELECT expr, expr, expr` is even **more** complicated-looking.

Let's look at details. All the syntax means there are a lot of places to add processing.

The From Clause

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

The FROM tables need to be iterable sequences.

```
list[dict[str, Any]]
```

```
from itertools import product
```

```
from_ = product(t1, t2, t3)
```

Yes. It's the Cartesian product.

Aha!

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

“Gotcha!”

- A *real* database doesn't do cartesian products all the time.
- It has fancy query algorithms and optimizations.

Aha!

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

“Gotcha!”

- A *real* database doesn't do cartesian products all the time.
- It has fancy query algorithms and optimizations.

“Your nonsense is clearly unworkable in general.”

Query optimization

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

- Requires extended syntax to suggest query optimizations.
- Require someone to design the right indexes.
- Requires detailed statistics on key distribution.

You **can** do this query optimization in Python, also.

Query optimization

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

- Requires extended syntax to suggest query optimizations.
- Require someone to design the right indexes.
- Requires detailed statistics on key distribution.

You **can** do this query optimization in Python, also.

We'll get to it.

The Where Condition

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

Context:

```
where = (row for row in from_ if c(row))
```

Or.

```
where = filter(c, from_)
```

```
def c(row: tuple[dict[str, Any], ...]) -> bool:
    t1, t2, t3 = row
    return (
        t1['rowid'] == t2['foreign_key']
        and t2['some_key'] == t3['whatever']
    )
```

The Select Clause

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

Context:

```
result = list(  
    (ex1(row_tuple), ex2(row_tuple), ...)   
    for row_tuple in where)
```

Or.

```
result = map(row_builder, where)
```

```
def e1(row: tuple[dict[str, Any], ...]) -> Any:  
    t1, t2, t3 = row  
    if t1['value'] % 2 == 0:  
        return t1['value'] // 2  
    else:  
        return t1['value'] * 3 + 1
```

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

Okay, From, Where, and Select clauses not awful.

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

Okay, From, Where, and Select clauses not awful.

The From clause needs work.

The Cartesian Product Problem

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

Databases do “cart-prod” joins all the time.

Very small tables are easier to fetch from disk into cache without wasting time on the additional index read.

Two common alternative algorithms:

- Sort-Merge Join
- Lookup Join

Sort-Merge Join

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

For very large tables.

- ① Sort each table into a consistent order by the join key for that table.
- ② Create row tuples for matching rows from each sorted table.

A variation on this can do any of the outer join algorithms.

See <https://toolz.readthedocs.io>; they offer `merge_sorted()`

Lookup Join

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

Create for many small tables and one big table.

The “Star Schema” design pattern.

Transform each small table into a Python dictionary.

Then, this:

```
from_ = (  
    (r, small_1[r['key_1']], small_2[r['key_2']])  
    for r in big_table  
)
```

Okay.

Fine.

The From, Where, and Select clauses have really fast pure-Python implementations.

What about Group by and Having?

Okay.

Fine.

The From, Where, and Select clauses have really fast pure-Python implementations.

What about Group by and Having?
They can't be simple.

Group By Clause(s)

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

The **Group By** process involves two separate things.

- The expressions (and column names) in the **GROUP BY** clause.

These expressions define a kind of `reduce()` function to build the groups.

- The Aggregate Functions in the **SELECT** clause then get applied to each group's subset of rows.

Syntax Oddity: Aggregates written in the **SELECT** clause.
Some in the **HAVING** clause.

Group By Implementation

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

- ① Partition into groups, `defaultdict(list)`.
For each row:
 - ① Create keys
 - ② Append row to groups
- ② For each group:
 - ① Evaluate all the aggregate functions.

This builds a new table from the keys and aggregate functions for each group.

Group By Implementation

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

```
from collections import defaultdict
from operator import itemgetter

# Partition

groups = defaultdict(list)
for row_tuple in where:
    key = (k_1(row_tuple), k_2(row_tuple), ...)
    groups[key].append(row_tuple)

# Aggregate

group_by = []
for key, group in groups:
    agg_1 = some_function(group)
    agg_2 = mean(row['value'] for row in group)
    agg_3 = sum(map(itemgetter('name'), group))
    group_by.append((key, agg_1, agg_2, agg_3))
```

Having Clause

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

The `Having` process is (nearly) the same as the `Where` process. It's an expression to filter the groups.

SQL syntax uses aggregate functions in the `Having` clause.

- These are yet more group-by aggregates.
- The result values are only used for filtering.

Conclusion

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

- SQL is helpful to summarize a desired result.
- SQL can be overused.
- A database is a lot of overhead. Avoid it.

Think of SQL as a design language.

Conclusion

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

- SQL is helpful to summarize a desired result.
- SQL can be overused.
- A database is a lot of overhead. Avoid it.

Think of SQL as a design language.
Not an implementation choice.

SQL

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

- SQL describes a pipeline of steps

From \rightarrow Where \rightarrow Select \rightarrow Group By \rightarrow Having

- Or. As nested functions

$$H\left(G_A\left(G_P\left(S\left(W\left(F(t_1, t_2, \dots)\right)\right)\right)\right)\right)$$

Nested functions seem confusing.

SQL to Python

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

In Python, Select is a stack of generator expressions

```
from_ = itertools.product(...)
where = filter(condition, from_)
select = map(row_builder, where)
groups = group_reduce(select)
aggregates = map(agg_row_builder, groups)
result = filter(having_condition, aggregates)
```

Most steps are lazy and don't compute big intermediate results.

SQL to Python

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

In Python, Select is a stack of generator expressions

```
from_ = itertools.product(...)
where = filter(condition, from_)
select = map(row_builder, where)
groups = group_reduce(select)
aggregates = map(agg_row_builder, groups)
result = filter(having_condition, aggregates)
```

Most steps are lazy and don't compute big intermediate results.
The `group_reduce()` function does compute a big result.

Call to Action

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

Stop using SQL as a data transformation tool.

Continue using SQL as a design aid.

More Information

Unlearning
SQL

S.Lott

SQL Overuse

SQL Design
Patterns

Python Imple-
mentation
Patterns

About Those
Join
Algorithms

Group By and
Having — The
Good Stuff

Conclusion

- *Unlearning SQL* (Available from Amazon and Lulu)
- <https://github.com/slott56/functional-SQL>
- <https://fosstodon.org/@slott56>