University of Trieste

*Data Management for Big Data* Course

Academic Year 2022–2023

# Data Warehouse case study

Davide Capone[*]     Sandro Junior Della Rovere[†]     Enrico Stefanel[‡]

Last update on 2023–06–22T14:30:15Z.

## 1   Introduction

The aim of this project is to study an efficient implementation of a suite of business oriented ad-hoc queries over the public TPC-H benchmark, which can be considered as a Big Data database, that has been implemented in Postgres.

### 1.1   TPC-H benchmark database

The TPC-H benchmark is a decision support benchmark that can be downloaded from the TPC official website. The data generator lets the user specify a *scale factor* in order to control the size of the resulted database. Our choices was to use a *scale factor* of 10, meaning that the overall database size is approximately 13 GB.

#### 1.1.1   Database statistics

The benchmark is composed by eight tables:

- CUSTOMER, with 16 columns and 1 500 000 tuples (312 MB);

- LINEITEM, with 32 columns and 59 986 052 tuples (11 GB); the main attributes that are going to be used are:

[*]davide.capone@studenti.units.it
[†]sandrojunior.dellarovere@studenti.units.it
[‡]enrico.stefanel@studenti.units.it

- $\circ$ `l_extendedprice` (1 351 462 distinct values, i.e. there is an average of 44 tuples with the same value, that range from 900.91 to 104 949.50),

- $\circ$ `l_discount` (11 distinct values, i.e. there is an average of 5 453 277 tuples with the same value, that range from 0.00 to 0.10),

- $\circ$ `l_returnflag` (which can assume values A→accepted, R→returned, N→not yet delivered; the percentage of tuples for A and R are almost 25 %, while the percentage of tuples where `l_returnflag` is N is about 50 %),

- $\circ$ `l_commitdate` (2466 distinct values, i.e. there is an average of 24 325 tuples with the same value, that range from 1992-01-31 to 1998-10-31),

- $\circ$ `l_receiptdate` (2555 distinct values, i.e. there is an average of 23 478 tuples with the same value, that range from 1992-01-03 to 1998-12-31);

- `NATION`, with 8 columns and 25 tuples (24 kB);

- `ORDERS`, with 18 columns and 1 500 000 tuples (2481 kB); the main attributes that are going to be used are:

  - $\circ$ `o_orderdate` (2406 distinct values, i.e. there is an average of 6234 tuples with the same value, that range from 1992-01-01 to 1998-08-02);

- `PART`, with 18 columns and 2 000 000 tuples (363 MB); the main attributes that are going to be used are:

  - $\circ$ `p_type` (150 distinct values, i.e. there is an average of 13 333 tuples with the same value);

- `PARTSUPP`, 10 columns and with 8 000 000 tuples (1535 MB);

- `REGION`, 6 columns and with 5 tuples (24 kB);

- `SUPPLIER`, 14 columns and with 100 000 tuples (20 MB).

Other attributes have been used, but statistics about them have been omitted for lack of usefulness (e.g., keys of the tables, for which the cardinality is exactly the cardinality of the corresponding table).

### 1.1.2  Database SQL definition

The SQL definition of the tables can be found on the official benchmark download.

# 2  Set of queries

## 2.1  Export/import revenue value

```
 1  WITH lineitem_orders AS (
 2      SELECT
 3          l_partkey,
 4          l_suppkey,
 5          o_orderdate,
 6          o_custkey,
 7          l_extendedprice,
 8          l_discount
 9      FROM lineitem JOIN orders ON (l_orderkey = o_orderkey)
10  ), customer_location AS (
11      SELECT
12          c_custkey,
13          c_name,
14          n_nationkey AS c_nationkey,
15          n_name AS c_nationname,
16          r_regionkey AS c_regionkey,
17          r_name AS c_regionname
18      FROM customer
19          JOIN nation ON (c_nationkey = n_nationkey)
20          JOIN region ON (n_regionkey = r_regionkey)
21  ), supplier_location AS (
22      SELECT
23          s_suppkey,
24          s_name,
25          n_nationkey AS s_nationkey,
26          n_name AS s_nationname,
27          r_regionkey AS s_regionkey,
28          r_name AS s_regionname
29      FROM supplier
30          JOIN nation ON (s_nationkey = n_nationkey)
31          JOIN region ON (n_regionkey = r_regionkey)
32  ), query1 AS (
33      SELECT
34          EXTRACT (YEAR FROM o_orderdate) AS _year,
35          EXTRACT (QUARTER FROM o_orderdate) AS _quarter,
36          EXTRACT (MONTH FROM o_orderdate) AS _month,
37          c_regionname,
38          c_nationname,
39          c_name,
40          s_regionname,
41          s_nationname,
42          s_name,
```

```
43            p_type,
44            SUM(l_extendedprice * (1 - l_discount)) AS revenue
45        FROM lineitem_orders
46            JOIN part ON l_partkey = p_partkey
47            JOIN supplier_location ON (s_suppkey = l_suppkey)
48            JOIN customer_location ON (c_custkey = o_custkey)
49        WHERE s_nationkey <> c_nationkey
50        GROUP BY
51            _year,
52            _quarter,
53            _month,
54            c_regionkey,
55            c_regionname,
56            c_nationkey,
57            c_nationname,
58            c_custkey,
59            c_name,
60            s_regionkey,
61            s_regionname,
62            s_nationkey,
63            s_nationname,
64            s_suppkey,
65            s_name,
66            p_type
67  )
68  SELECT * FROM query1;
```

## 2.2 Late delivery

It is asked to retrieve the number of orders where at least one "lineitem" has been received later than the committed date. The aggregation should be performed with the Month → Year roll-up, and the (Customer's) Nation → Region roll-up.

```
1   WITH lineitem_orders AS (
2       SELECT
3           o_orderkey,
4           l_partkey,
5           l_suppkey,
6           o_orderdate,
7           o_custkey,
8           l_commitdate,
9           l_receiptdate
10      FROM lineitem JOIN orders ON (l_orderkey = o_orderkey)
```

```
11  ), customer_location AS (
12      SELECT
13          c_custkey,
14          n_nationkey AS c_nationkey,
15          n_name AS c_nationname,
16          r_regionkey AS c_regionkey,
17          r_name AS c_regionname
18      FROM customer
19          JOIN nation ON (c_nationkey = n_nationkey)
20          JOIN region ON (n_regionkey = r_regionkey)
21  ), query2 AS (
22  SELECT
23      EXTRACT(YEAR FROM o_orderdate) AS _year,
24      EXTRACT(MONTH FROM o_orderdate) AS _month,
25      c_regionname,
26      c_nationname,
27      COUNT(DISTINCT(o_orderkey)) AS orders_no
28  FROM lineitem_orders
29      JOIN part ON l_partkey = p_partkey
30      JOIN customer_location ON (c_custkey = o_custkey)
31  WHERE
32      l_receiptdate > l_commitdate
33      -- AND _month = 1
34      -- AND p_type = 'PROMO BURNISHED COPPER'
35  GROUP BY
36      _year,
37      _month,
38      c_regionkey,
39      c_regionname,
40      c_nationkey,
41      c_nationname
42  )
43  SELECT * FROM query2;
```

## 2.3   Returned item loss

It is asked to retrieve the *revenue loss* for customers who might be having problems with the parts that are shipped to them, where a *revenue loss* is defined as

$$SUM(l\_extendedprice*(1-l\_discount))$$

for all qualifying *lineitems*.

```
1  WITH lineitem_orders AS (
```

```
 2      SELECT
 3          o_orderkey,
 4          o_orderdate,
 5          o_custkey,
 6          l_extendedprice,
 7          l_discount,
 8          l_returnflag
 9      FROM lineitem JOIN orders ON (l_orderkey=o_orderkey)
10  ),
11  query3 AS (
12  SELECT
13      EXTRACT(YEAR FROM o_orderdate) AS _year,
14      EXTRACT(QUARTER FROM o_orderdate) AS _quarter,
15      EXTRACT(MONTH FROM o_orderdate) AS _month,
16      c_name,
17      SUM(l_extendedprice*(1-l_discount)) AS returnloss
18  FROM
19      lineitem_orders
20      JOIN customer ON (o_custkey=c_custkey)
21  WHERE
22      l_returnflag='R'
23      -- AND c_name='Customer#000129976'
24      -- AND EXTRACT(QUARTER FROM o_orderdate) = 1
25  GROUP BY
26      _year,
27      _quarter,
28      _month,
29      c_custkey,
30      c_name
31  )
32  SELECT * FROM query3;
```

Five independent runs of the above query obtained the following execution times: 753 541.744 ms, 672 530.120 ms, 624 276.525 ms, 615 741.447 ms and 634 262.713 ms.

# 3   Indexes design

. . .

# 4   Materialisation

. . .

# 5 Conclusions

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Pellentesque mollis odio felis, eu aliquam velit congue at. Morbi et aliquet risus. Sed varius, diam vel semper varius, nibh mi bibendum est, vel molestie magna elit nec elit. Morbi ut facilisis mi. Suspendisse tristique porta purus, ac convallis tellus porta et. Maecenas ac semper justo, ut tincidunt turpis. Proin euismod viverra nisi iaculis accumsan. Suspendisse hendrerit suscipit purus sed viverra. Nullam gravida suscipit sagittis. Donec at pulvinar nunc. Donec consequat nibh elit, ornare posuere nisi feugiat laoreet. Vivamus nec imperdiet justo.
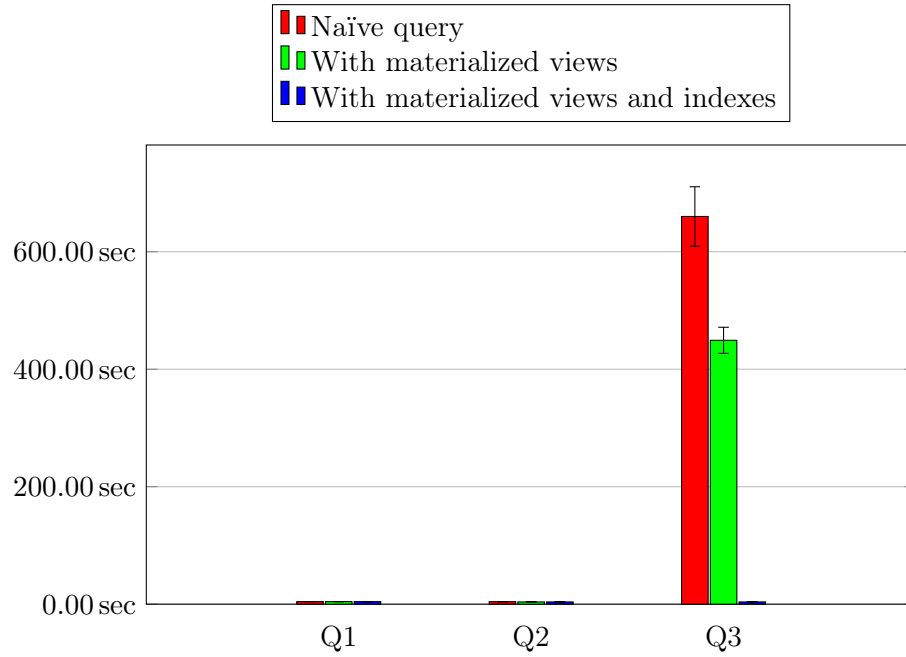


Figure 1: Query timings

Mauris vitae eros lacus. Etiam facilisis orci ac nunc tincidunt, sed vehicula sem ultrices. Sed dignissim velit elit, vel dictum tortor feugiat eget. Nulla nec sagittis quam. Cras diam mauris, suscipit vel urna in, congue laoreet diam. In consequat, nisi non iaculis interdum, risus lacus auctor nisi, in sagittis massa libero in nibh. Phasellus venenatis urna eu nibh maximus ullamcorper. Pellentesque tincidunt, nulla vel ultrices malesuada, justo velit sollicitudin orci, vitae euismod ipsum nisl non magna.