

Scalable Data Engineering, Exercise 3

HENRY HAUSTEIN

Task 1

- (a) True
- (b) False, maintenance is different
- (c) False, it is an NP-hard problem
- (d) False, it depends on the use case
- (e) True

Task 2

- (a) SQL:

```
1 EXPLAIN(SELECT * FROM ORDERS)
```

- (b) SQL:

```
1 SELECT * FROM ORDERS WHERE O_CLERK = "Clerk#000000322"
```

took 600ms.

- (c) SQL:

```
1 CREATE INDEX ocleak_idx ON ORDERS(O_CLERK)
```

- (d) Query from (b) then takes around 300ms

Task 3

- (a) SQL:

```
1 CREATE MATERIALIZED VIEW custperreg AS
2   SELECT REGION.R_NAME, COUNT(*)
3   FROM CUSTOMER, NATION, REGION
4   WHERE
5     CUSTOMER.C_NATIONKEY = NATION.N_NATIONKEY AND
6     NATION.N_REGIONKEY = REGION.R_REGIONKEY
7   GROUP BY REGION.R_NAME
8   WITH DATA
```

(b) SQL:

```
1  INSERT INTO REGION(R_REGIONKEY, R_NAME, R_COMMENT) VALUES (5,
    "AUSTRALIA", "down under")
2
3  REFRESH MATERIALIZED VIEW custperreg
```

(c) In the materialized view Australia is missing. That's why there are no customers for this region and since we are doing inner joins for the materialized view there are no joining partners for Australia. To fix this:

```
1  DROP MATERIALIZED VIEW custperreg
2
3  CREATE MATERIALIZED VIEW custperreg AS
4      SELECT REGION.R_NAME, COUNT(CUSTOMER.C_CUSTKEY)
5      FROM
6          CUSTOMER RIGHT OUTER JOIN (
7              NATION RIGHT OUTER JOIN REGION
8              ON(NATION.N_REGIONKEY = REGION.R_REGIONKEY)
9          )
10         ON(CUSTOMER.C_NATIONKEY = NATION.N_NATIONKEY)
11     GROUP BY REGION.R_NAME
12     WITH DATA
```

Task 4

(a) We create a materialized view from the subquery:

```
1  CREATE MATERIALIZED VIEW linepart AS
2      SELECT
3          PART.P_NAME,
4          SUM(LINEITEM.L_QUANTITY) AS qty
5      FROM PART, LINEITEM
6      WHERE PART.P_NAME = LINEITEM.L_PARTKEY
7      GROUP BY PART.P_NAME
8      WITH DATA
9
10     SELECT PART.P_NAME, qty, RANK() OVER(ORDER BY qty DESC)
11     FROM linepart
```

Creating the materialized view takes some time but then the actual query is done in about a second (from 10 seconds before).

(b) SQL:

```
1  CREATE MATERIALIZED VIEW linepart2 AS
2      SELECT
3          PART.P_NAME,
4          SUM(LINEITEM.L_QUANTITY) AS qty,
5          LINEITEM.L_SUPPKEY
6      FROM PART, LINEITEM
7      WHERE PART.P_NAME = LINEITEM.L_PARTKEY
8      GROUP BY PART.P_NAME, LINEITEM.L_SUPPKEY
```

9 WITH DATA

The query from (a) becomes then

```
1 SELECT PART.P_NAME, SUM(qty), RANK() OVER(ORDER BY SUM(qty)
   DESC)
2 FROM linepart2
```

with roughly 4 seconds run time. The other query becomes

```
1 SELECT
2     PART.P_NAME ,
3     NATION.N_NAME ,
4     SUM(qty),
5     RANK() OVER(
6         PARTITION BY PART.P_NAME
7         ORDER BY SUM(qty) DESC
8     )
9 FROM linepart2, NATION, SUPPLIER
10 WHERE
11     LINEITEM.L_SUPPKEY = SUPPLIER.S_SUPPKEY AND
12     SUPPLIER.S_NATIONKEY = NATION.N_NATIONKEY
13 GROUP BY PART.P_NAME, NATION.N_NAME
```

with now 21 seconds runtime (was 60 seconds before).