

Name

CWID

Final Exam

May 4th, 2020

8:00-11:00

CS520 - Data Integration, Warehousing, and Provenance

Results

Please leave this empty!

1.1

1.2

Sum

Instructions

- **When writing a query, write the query in a way that it would work over all possible database instances and not just for the given example instance!**

Consider the following datawarehouse schema (star schema) and partial example instance. There is a single fact table (**sales**) about sales of items. Each row in this fact table stores the quantity of a certain product (e.g., 3 Samson Galaxy phones) sold at a particular location and time to a particular customer. There are four dimension tables corresponding to the following dimensions:

- **Time** with three levels (year, month, day)
- **Location** with four levels (state, city, zip, street)
- **Customer** with one level (name). Customers are associated with a location through a foreign key (attribute CLID) to dimension **Location**.
- **Product** with three levels (category, brand, pname, price) where pname is the finest granularity and brand and category are not comparable (some brands can have products from multiple categories and categories obviously can contain have products from different brands). The same holds for price and brand and price and category.

sales

TID	LID	CID	PID	numItems
1	4	1	1	15
2	1	5	2	10
100	1	76	4	22
...

timeDim

TID	year	month	day
1	2010	1	1
2	2010	1	2
...
...	2018	5	1

locationDim

LID	state	city	zip	street
1	Illinois	Chicago	60616	10 W 31st
2	Illinois	Chicago	60615	900 Cottage Grove
3	Lousiana	New Orleans	42345	12 Mark street
...

customerDim

CID	cname	CLID
1	Noekig	1
2	Prokig	3
...

productDim

PID	category	brand	pname	price
1	computers	Apple	MacBook	1300
2	computers	Dell	Inspire	1000
3	smartphones	Samsung	Galaxy 1	600
...

Hints:

- Attributes with black background form the primary key of a relation (e.g., PID for relation `productDim`)
- Attributes LID, TID, PID, and CID in the fact table are foreign keys to the dimension tables
- Attribute CLID of relation `customerDim` is a foreign key to relation `locationDim`

Part 1.1 Data Warehousing (Total: 60 Points)

Recall that you should write all queries according to the schema and not according to the example instance.

Question 1.1.1 (8 Points)

Write an SQL query that returns for each city, the 3 zip codes with the highest total revenue (calculated as the sum over the number of items sold multiplied by their price).

Solution

```
WITH zipRev AS (  
    SELECT sum(numItems * price) AS rev, city, zip  
    FROM sales s, locationDim l, productDim p  
    WHERE s.LID = l.LID AND s.PID = p.PID  
    GROUP BY city, zip)  
SELECT city, zip  
FROM (SELECT *, row_number() OVER (PARTITION BY city ORDER BY rev DESC) AS pos  
      FROM zipRev) sub  
WHERE pos <= 3;
```

Question 1.1.2 (8 Points)

Write an SQL query that returns for each product category and year, the accumulative sum of the revenue over the months, i.e., the result for Feb 2010 would be the sum of revenue for 2010 Jan and Feb, for Mar 2010 return the sum of 2010 Jan, Feb, Mar, and so on.

Solution

```
SELECT category , year , month , sum(rev) OVER (PARTITION BY category , year
                                                ORDER BY month ASC) AS accrev
FROM (SELECT sum(numItems * price) AS rev , category , year , month
      FROM sales s , productDim p , timeDim t
      WHERE s.PID = p.PID AND s.TID = t.TID
      GROUP BY category , year , month) mrev
```

Question 1.1.3 (8 Points)

Write an SQL query that returns the total revenue for each level of the time dimension.

Solution

```
SELECT sum(numItems * price) as rev ,
       year ,
       month ,
       day ,
       GROUPING(year) AS grpYear ,
       GROUPING(month) AS grpMonth ,
       GROUPING(day) AS grpDay
FROM sales s, timeDim t, productDim p
WHERE s.TID = t.TID AND s.PID = p.PID
GROUP BY ROLLUP(year ,month ,day);
```

Question 1.1.4 (9 Points)

Write an SQL query that returns the three customers per state (location of the customer) that generated the largest revenue in their state. For example, for the result for Illinois only include customers from Illinois and only consider the revenue that they produced (by buying items) in Illinois.

Solution

```
WITH custStateRev AS (  
    SELECT sum(numItems * price) AS rev, cname, state  
    FROM sales s, customerDim c, locationDim l, productDim p  
    WHERE s.CID = c.CID AND s.LID = l.LID AND c.CLID = l.LID AND s.PID = p.PID  
    GROUP BY cname, state  
)  
SELECT state, cname, rev  
FROM (SELECT *, row_number() OVER (PARTITION BY state ORDER BY rev DESC) AS pos  
      FROM custStateRev) sub  
WHERE pos <= 3;
```

Question 1.1.5 (9 Points)

Return the three years with the largest growth factor. The growth factor for a year should be calculated as the sum of the revenue for that year divided by the average of the sums of revenue for the three preceding years.

Solution

```
WITH yearlyRev AS (  
    SELECT sum(numItems * price) AS rev, year  
    FROM sales s, timeDim t, productDim p  
    WHERE s.TID = t.TID AND s.PID = p.PID  
    GROUP BY year),  
yearGrowth AS (  
    SELECT year, rev / avg(rev) OVER (ORDER BY year ASC  
                                     ROWS BETWEEN 3 PRECEDING AND 1 PRECEDING) AS gf  
    FROM yearlyRev  
)  
SELECT year, gf  
FROM yearGrowth  
ORDER BY gf DESC  
LIMIT 3;
```


Question 1.1.6 (9 Points)

Write an SQL query that returns for brands Apple, Dell, and IBM the number of years where items of this brand sold during the year produced the greatest revenue among all brands.

Solution

```
WITH yearlyBrandRev AS (  
    SELECT sum(numItems * price) AS rev, brand, year  
    FROM sales s, productDim p, timeDim t  
    WHERE s.PID = p.PID AND s.TID = t.TID  
    GROUP BY brand, year  
),  
bestPerYear AS (  
    SELECT max(rev) AS maxrev, year  
    FROM yearlyBrandRev  
    GROUP BY year  
)  
SELECT brand, sum(CASE WHEN rev = maxrev THEN 1 ELSE 0 END) AS numYearsBest  
FROM yearlyBrandRev b, bestPerYear m  
WHERE brand IN ('Apple', 'Dell', 'IBM') AND b.year = m.year  
GROUP BY brand;
```

Question 1.1.7 (9 Points)

Write an SQL query that returns for each customer the percentage of items they brought per brand. For example, if a customer buys 10 Apple products and 50 products in total then you should return 20% as the fraction for Apple for this customer.

Solution

```
WITH custBrand AS (  
    SELECT sum(numItems) AS brandItems, cname, brand  
    FROM sales s, productDim p, customerDim c  
    GROUP BY cname, brand  
)  
SELECT (brandItems / ttlItems) * 100.0 AS perc, cname, brand  
FROM (SELECT brandItems,  
            sum(brandItems) OVER (PARTITION BY cname) AS ttlItems,  
            cname,  
            brand  
FROM custBrand) AS cttl;
```

Part 1.2 Provenance (Total: 40 Points)

For each of the queries shown in the following compute the provenance of all of their result tuples produced over the database shown below. Calculate provenance for these provenance models:

- Why-Provenance
- Minimal Why-Provenance
- Provenance Polynomials

Before showing the provenance, first write down the results of the query and label the result tuples t_1, t_2, \dots, t_n .

Consider the following database schema and instance:

hotel

hname	location	owner	roomsnt	
Astor	Chicago	Bob	20	h_1
Blackstone	Chicago	Bob	120	h_2
Seaside	Miami	Alice	12	h_3

booking

hotel	tourist	startdate	enddate	rooms	rate	
Astor	Peter	2004-01-01	2004-01-03	1	230	b_1
Astor	Tilda	2004-01-05	2004-01-06	1	200	b_2
Seaside	Peter	2004-02-20	2004-02-30	2	75	b_3
Blackstone	Alice	2004-03-01	2004-03-05	1	140	b_4

person

name	age	location	
Peter	43	Chicago	p_1
Bob	24	Madison	p_2
Alice	25	Chicago	p_3
Tilda	32	Miami	p_4

location

loc	state	country	
Chicago	43	USA	l_1
Madison	24	USA	l_2
Miami	25	USA	l_3

Question 1.2.1 (7 Points)

$$\pi_{location}(\sigma_{roomcnt > 15}(hotel))$$

Solution

Result relation:

location
Chicago

 t_1

Why provenance:

location
Chicago

 $\{\{h_1\}, \{h_2\}\}$

Minimal Why provenance:

location
Chicago

 $\{\{h_1\}, \{h_2\}\}$

Provenance Polynomials:

location
Chicago

 $h_1 + h_2$

Question 1.2.2 (9 Points)

$$q_1 \stackrel{def}{=} booking \bowtie_{hotel=hname} hotel \bowtie_{loc=location} location$$

$$q_2 \stackrel{def}{=} \rho_{plocation \leftarrow location}(person) \bowtie_{plocation=ploc} \rho_{ploc \leftarrow loc, pstate \leftarrow state, pcountry \leftarrow country}(location)$$

$$q \stackrel{def}{=} \pi_{hname, plocation, pcountry, location, country}(q_1 \bowtie_{tourist=name} q_2)$$

Solution

Result relation:

name	plocation	pcountry	location	country	
Astor	Chicago	USA	Chicago	USA	t_1
Astor	Miami	USA	Chicago	USA	t_2
Blackstone	Chicago	USA	Chicago	USA	t_3
Seaside	Chicago	USA	Miami	USA	t_4

Why provenance:

name	plocation	pcountry	location	country	
Astor	Chicago	USA	Chicago	USA	$\{\{b_1, h_1, p_1, l_1\}\}$
Astor	Miami	USA	Chicago	USA	$\{\{b_2, h_1, p_4, l_1, l_3\}\}$
Blackstone	Chicago	USA	Chicago	USA	$\{\{b_4, h_2, p_3, l_1\}\}$
Seaside	Chicago	USA	Miami	USA	$\{\{b_3, h_3, p_1, l_1, l_3\}\}$

Minimal Why provenance:

name	plocation	pcountry	location	country	
Astor	Chicago	USA	Chicago	USA	$\{\{b_1, h_1, p_1, l_1\}\}$
Astor	Miami	USA	Chicago	USA	$\{\{b_2, h_1, p_4, l_1, l_3\}\}$
Blackstone	Chicago	USA	Chicago	USA	$\{\{b_4, h_2, p_3, l_1\}\}$
Seaside	Chicago	USA	Miami	USA	$\{\{b_3, h_3, p_1, l_1, l_3\}\}$

Provenance Polynomials:

name	plocation	pcountry	location	country	
Astor	Chicago	USA	Chicago	USA	$b_1 \cdot h_1 \cdot p_1 \cdot l_1^2$
Astor	Miami	USA	Chicago	USA	$b_2 \cdot h_1 \cdot p_4 \cdot l_1 \cdot l_3$
Blackstone	Chicago	USA	Chicago	USA	$b_4 \cdot h_2 \cdot p_3 \cdot l_1^2$
Seaside	Chicago	USA	Miami	USA	$b_3 \cdot h_3 \cdot p_1 \cdot l_1 \cdot l_3$

Question 1.2.3 (8 Points)

$$q \stackrel{def}{=} \pi_{location}(\sigma_{country=USA}(person \bowtie_{location=loc} location))$$

Solution

Result relation:

location	
Chicago	t_1
Madison	t_2
Miami	t_3

Why provenance:

location	
Chicago	$\{\{p_1, l_1\}, \{p_3, l_1\}\}$
Madison	$\{\{p_2, l_2\}\}$
Miami	$\{\{p_4, l_3\}\}$

Minimal Why provenance:

location	
Chicago	$\{\{p_1, l_1\}, \{p_3, l_1\}\}$
Madison	$\{\{p_2, l_2\}\}$
Miami	$\{\{p_4, l_3\}\}$

Provenance Polynomials:

location	
Chicago	$p_1 \cdot l_1 + p_3 \cdot l_1$
Madison	$p_2 \cdot l_2$
Miami	$p_4 \cdot l_3$

Question 1.2.4 (8 Points)

$$q_1 \stackrel{def}{=} \pi_{hotel, date}(\rho_{date \leftarrow startdate}(booking))$$

$$q_2 \stackrel{def}{=} \pi_{hotel, date}(\rho_{date \leftarrow enddate}(booking))$$

$$q \stackrel{def}{=} q_1 \cup q_2$$

Solution

Result relation:

hotel	date	
Astor	2004-01-01	t_1
Astor	2004-01-05	t_2
Astor	2004-01-03	t_3
Astor	2004-01-06	t_4
Seaside	2004-02-20	t_5
Seaside	2004-02-30	t_6
Blackstone	2004-03-01	t_7
Blackstone	2004-03-05	t_8

Why provenance:

hotel	date	
Astor	2004-01-01	$\{\{b_1\}\}$
Astor	2004-01-03	$\{\{b_1\}\}$
Astor	2004-01-05	$\{\{b_2\}\}$
Astor	2004-01-06	$\{\{b_2\}\}$
Seaside	2004-02-20	$\{\{b_3\}\}$
Seaside	2004-02-30	$\{\{b_3\}\}$
Blackstone	2004-03-01	$\{\{b_4\}\}$
Blackstone	2004-03-05	$\{\{b_4\}\}$

Minimal Why provenance:

hotel	date	
Astor	2004-01-01	$\{\{b_1\}\}$
Astor	2004-01-03	$\{\{b_1\}\}$
Astor	2004-01-05	$\{\{b_2\}\}$
Astor	2004-01-06	$\{\{b_2\}\}$
Seaside	2004-02-20	$\{\{b_3\}\}$
Seaside	2004-02-30	$\{\{b_3\}\}$
Blackstone	2004-03-01	$\{\{b_4\}\}$
Blackstone	2004-03-05	$\{\{b_4\}\}$

Provenance Polynomials:

hotel	date	
Astor	2004-01-01	b_1
Astor	2004-01-03	b_1
Astor	2004-01-05	b_2
Astor	2004-01-06	b_2
Seaside	2004-02-20	b_3
Seaside	2004-02-30	b_3
Blackstone	2004-03-01	b_4
Blackstone	2004-03-05	b_4

Question 1.2.5 (8 Points)

$$q \stackrel{def}{=} \pi_{hname, name}(person \bowtie hotel)$$

Solution

Result relation:

hname	name	
Astor	Peter	t_1
Astor	Alice	t_2
Blackstone	Peter	t_3
Blackstone	Alice	t_4
Seaside	Tilda	t_5

Why provenance:

hname	name	
Astor	Peter	$\{\{h_1, p_1\}\}$
Astor	Alice	$\{\{h_1, p_3\}\}$
Blackstone	Peter	$\{\{h_2, p_1\}\}$
Blackstone	Alice	$\{\{h_2, p_3\}\}$
Seaside	Tilda	$\{\{h_3, p_4\}\}$

Minimal Why provenance:

hname	name	
Astor	Peter	$\{\{h_1, p_1\}\}$
Astor	Alice	$\{\{h_1, p_3\}\}$
Blackstone	Peter	$\{\{h_2, p_1\}\}$
Blackstone	Alice	$\{\{h_2, p_3\}\}$
Seaside	Tilda	$\{\{h_3, p_4\}\}$

Provenance Polynomials:

hname	name	
Astor	Peter	$h_1 \cdot p_1$
Astor	Alice	$h_1 \cdot p_3$
Blackstone	Peter	$h_2 \cdot p_1$
Blackstone	Alice	$h_2 \cdot p_3$
Seaside	Tilda	$h_3 \cdot p_4$

