

Homework 2

*Handed Out: September 12, 2017**Due: September 26, 2017 11:59 pm*

1 General Instructions

- This assignment is due at 11:59 PM on the due date. We will be using Sakai (<https://sakailogin.nd.edu/portal/site/FA17-CSE-40647-CX-01>) for collecting this assignment. Contact TA if you face technical difficulties in submitting the assignment. We shall NOT accept any late submission!
- The homework MUST be submitted in pdf format. Handwritten answers are not acceptable. Name your pdf file as YourNetid-HW2.pdf
- You need to explain the logic of your answer/result for every question. A result/answer without any explanation will not receive any point.
- It is OK to discuss the problems with the TA and your classmates, however, it is NOT OK to work together or share code. Plagiarism is an academic violation to copy, to include text from other sources, including online sources, without proper citation. To get a better idea of what constitutes plagiarism, consult the Honor code on academic integrity violations, including examples, and recommended penalties. There is a zero tolerance policy on academic integrity violations. Any student found to be violating this code will be subject to disciplinary action.
- Please use Piazza if you have questions about the homework. Also feel free to send TA emails and come to office hours.

2 Question 1 (25 points)

Suppose the base cuboid of a data cube contains two cells

$$(a_1, a_2, a_3, a_4, \dots, a_{10}) : 1, (a_1, b_2, a_3, b_4, \dots, b_{10}) : 1$$

where $a_i \neq b_i$ for any i . Obviously here we have 10 dimensions and each has only one level (no concept hierarchy).

1. How many **nonempty cuboids** are there in this data cube?
2. How many **nonempty aggregate closed cells** are there in this data cube?
3. How many **nonempty aggregate cells** are there in this data cube?
4. If we set minimum support = 2, how many **nonempty aggregate cells** are there in the corresponding iceberg cube?

Solution:

Treated second base cell as $(a_1, b_2, a_3, b_4, a_5, b_6, a_7, b_8, a_9, b_{10}) : 1$

1. (6') **Answer: 2^{10} .** Since we have 10 dimensions with no concept hierarchy, there are 2^{10} cuboids and all of them should not be empty.
2. (6') **Answer: 1.** There are 3 closed cells, including the two base cells and $(a_1, *, a_3, *, a_5, *, a_7, *, a_9, *)$. But only the latter one is an **aggregate** closed cell.
3. (6') **Answer: 2014.** For each base cell, there are $2^{10} - 1$ aggregated cells. However, there are 2^5 cells that are counted twice since there are 5 common dimensions. Therefore, the total number of nonempty aggregate cells is $2 \cdot (2^{10} - 1) - 2^5 = 2014$.
4. (6') **Answer: 2^5 .** These two base cells have common value in 5 dimensions; therefore, there are 2^5 nonempty cells with support = 2 and all of them are aggregate cells.

Treated second base cell as $(a_1, b_2, a_3, b_4, b_5, b_6, b_7, b_8, b_9, b_{10}) : 1$

1. (6') **Answer: 2^{10} .** Since we have 10 dimensions with no concept hierarchy, there are 2^{10} cuboids and all of them should not be empty.
2. (6') **Answer: 1.** There are 3 closed cells, including the two base cells and $(a_1, *, a_3, *, *, *, *, *, *, *)$. But only the latter one is an **aggregate** closed cell.
3. (6') **Answer: 2042.** For each base cell, there are $2^{10} - 1$ aggregated cells. However, there are 2^2 cells that are counted twice since there are 2 common dimensions. Therefore, the total number of nonempty aggregate cells is $2 \cdot (2^{10} - 1) - 2^2 = 2042$.
4. (6') **Answer: 4.** These two base cells have common value in 2 dimensions; therefore, there are 2^2 nonempty cells with support = 2 and all of them are aggregate cells.

3 Question 2 (25 points)

Which of the following algorithms: (i) Multiway array aggregation, (ii) BUC, cannot support the following operations efficiently? and explain why. This could be multiple-choice questions.

1. Computing an iceberg cube;
2. Supporting efficient OLAP query processing on a large dataset with 50 dimensions.

Solution:

1. (12') **Answer:** Multiway Array aggregation cannot support iceberg computation since it computes aggregation bottom-up and Apriori principle/pruning cannot be applied here. (Note that shell-fragment can support)
2. (13') **Answer:** Multiway Array and BUC – 50 dimensions is too many for both algorithms to support.

4 Question 3 (25 points)

Assume a base cuboid of 10 dimensions contains only three base cells¹:

- $(a_1, d_2, d_3, d_4, \dots, d_{10}) : 1,$
- $(d_1, b_2, d_3, d_4, \dots, d_{10}) : 1,$
- $(d_1, d_2, c_3, d_4, \dots, d_{10}) : 1,$

where $a_1 \neq d_1$, $b_2 \neq d_2$ and $c_3 \neq d_3$. The measure of the cube is *count*. Here we have 10 dimensions and each has only one level (no concept hierarchy).

1. How many **nonempty cuboids** will a full data cube contain?
2. How many **nonempty aggregate cells** will a full cube contain?
3. How many **nonempty aggregate cells** will an iceberg cube contain with the condition $count \geq 2$?

Solution:

1. (8') **Answer:** 2^{10}
2. (8') **Answer:** (i) Each cell generates $2^{10} - 1$ nonempty aggregated cells, thus in total we should have $3 * 2^{10} - 3$ cells with overlaps removed. (ii) We have $3 * 2^7$ cells overlapped once (thus count 2) and $1 * 2^7$ (which is $(*, *, *, d_4, \dots, d_{10})$) overlapped twice (thus count 3). Thus we should remove in total $1 * 3 * 2^7 + 2 * 1 * 2^7 = 5 * 2^7$ overlapped cells. (iii) Thus we have: $3 * 8 * 2^7 - 5 * 2^7 - 3 = 19 * 2^7 - 3$.
3. (9') **Answer:** (i) $(*, *, d_3, d_4, \dots, d_9, d_{10})$ has count 2 since it is generated by both cell 1 and cell 2; similarly, we have (ii) $(*, d_2, *, d_4, \dots, d_9, d_{10}) : 2$, (iii) $(*, *, d_3, d_4, \dots, d_9, d_{10}) : 2$; and (iv) $(*, *, *, d_4, \dots, d_9, d_{10}) : 3$. Therefore, we have $4 * 2^7 = 2^9$.

5 Question 4 (25 points)

We have a data array containing 3 dimensions A, B and C shown in Figure 1. The 3-D array is divided into 27 small chunks. Each dimension is divided into 3 equally sized partitions. The cardinality (size) of the dimensions A, B, and C is 900, 300, and 600. Since we divide each dimension into 3 parts with equal size, the sizes of the chunks on dimensions A, B, and C are 300, 100, and 200 respectively. Suppose we want to use **Multiway Array Aggregation Computation** to materialize the 2-D cuboids AB, AC and BC.

1. What is the ordering of chunk scanning that achieves the maximum computation efficiency, i.e. requires the least memory units?

¹Here when we say "only three base cells", we assume that we have ONE transaction in each of the base cells AND we only have these THREE transactions in the cuboid.

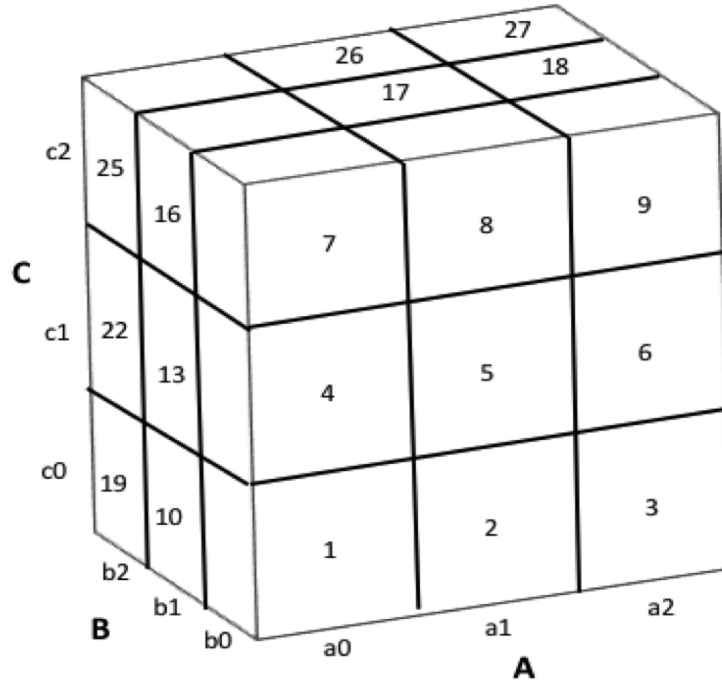


Figure 1: 3-D array of a data cube containing dimensions A, B, C and is divided into 27 small chunks.

- Following the ordering you give in the above subquestion, what is the minimum memory requirement for holding all the 2-D planes?

Solution:

- (12') **Answer:** 1, 10, 19, 4, 13, 22, 7, 16, 25, 2, 11, 20, 5, 14, 23, 8, 17, 26, 3, 12, 21, 6, 15, 24, 9, 18, 27.
Any correct answer will have full points, even like 8, 26, 17 (now aggregate B dimension for a1,c2); 11, 20, 2; 14, 23, 5; (now aggregate C dimension for a1,b0-b1-b2); 3, 12, 21; 6, 15, 24; 9, 18, 27; 1, 10, 19; 4, 13, 22; 7, 16, 25.
- (13') **Answer:** $300 \times 200 (AC) + 300 \times 300 (AB) + 300 \times 600 (BC) = 330000$.