

BUC

Olione Lonie

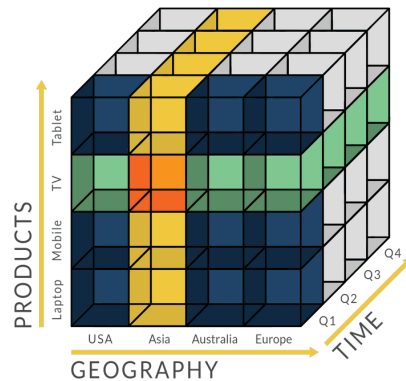
December 2022

1 Introduction

Les data-cubes sont utilisés pour présenter les données massives et les traiter plus facilement .

Un cube se partage en sous-cube selon n dimensions , ces dimensions sont constituées par un ensemble d'éléments , par exemple nous pourrions avoir la dimension produite qui est constitués de TV , brosse à dents , etc .

Dans ce rapport, nous étudierons la méthode **BUC** qui fut proposée par Kevin Beyer et Raghu Ramakrishnan dans : [Bottom-Up Computation of Sparse and Iceberg CUBES](#)



1.1 Intérêts :

Les entreprises ont aujourd'hui grand besoin d'avoir des méthodes rapides et efficaces pour répondre aux besoins de leur client .

En effet il a été montré dans le [Rapport Universel de l'Oreal 2021](#) que les clients demandent à leur marque de proposer des produits qui leur sont adaptés .

De plus d'après ce rapport les tendances des produits mais aussi les ralentissements de la chaîne de production pour chaque gamme doivent être déduites des reportings de chaque site . Les méthodes classiques d'extraction de connaissance (**ETL**) étant trop coûteuses pour traiter une si grande quantité de données de nouvelles méthodes doivent être développées .

C'est dans cet environnement que se place l'**Algorithme BUC** .

1.2 Principe :

Fonctions SQL :

Classiquement, les données sont découpées (*slicing*) grâce à la fonction *GROUP BY*. *GROUP BY* permet de regrouper les éléments partageant une même valeur d'attribut.

Pour éviter d'avoir à traiter l'entièreté de la base de données, on peut adjoindre une condition à ce *slicing* grâce au mot-clé *HAVING*.

Ceci jette les bases de l'approche iceberg-CUBE : plutôt que de créer l'ensemble du CUBE (agrégation de toutes les données pour l'ensemble des dimensions) on agrégera les données si et seulement si elle vérifie une condition fixée à l'avance, par exemple un *threshold*.

La Threshold Condition :

Un *threshold* est une borne qui délimite les informations pertinentes des informations peu utiles.

Imaginons qu'un commerce souhaite savoir quels sont les produits tendance pour chacune des saisons, alors il doit additionner (agrégation) l'ensemble de ses ventes pour chaque produit (*GROUP BY* product) et chaque saison (*GROUP BY* season) puis de calculer l'ensemble des indicateurs qu'il souhaite calculer (*SELECT {SUM, MAX, AVG, ...}*).

La méthode **BUC** assume qu'il est inutile de calculer les indicateurs pour les produits "peu vendus", car ils n'influencent pas la tendance globale, ceci justifie l'emploi d'un *threshold* :

$$HAVING SUM(*) > threshold$$

Ce constat permet de diminuer drastiquement le temps de calcul des indicateurs car seules les catégories/produits les plus importants seront calculés.

2 Concept :

2.1 Outils utilisés :

Comme nous l'avons vu la méthode **BUC** s'appuie en grande partie sur le "*thresholding*" il sera appelé support minimal (abrégié en *min - sup*).

Nous utiliserons la mère de la fonction *GROUP BY* : la fonction *CUBE*, elle calcule l'ensemble des *GROUP BY* sur les n dimensions de la données.

Par exemple, un enregistrement avec les dimensions produit, saison et localisation se verrait agréger par *CUBE* sur les dimensions (produit), (produit, saison), (produit, saison, localisation), (saison), (localisation) et (localisation, saison).

Ceci correspond à l'ensemble des permutations de (produit, saison, localisation).

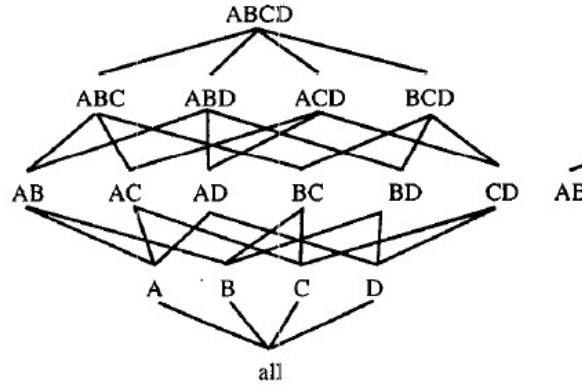
Un *CUBE* soumis à un *thresholding* sera appelé un iceberg-CUBE car seule la partie émergée (importante) des données se fera analysée en détail.

2.2 Methode :

2.2.1 Arbre :

Même avec ce *thresholding* l'opération *CUBE* reste très gourmande , c'est pourquoi une autre de nos préoccupations sera de la partitionner en plusieurs sous-problème résolvable par des machines indépendantes (pipeline cf cloud computing) .

Ce partitionnement va s'appuyer sur le principe de l'algèbre relationnel , rappelez vous un *CUBE* se divise en n dimension par exemple : **A** , **B** , **C** et **D** .



Si on part de la racine **ABCD** alors le graphe se lit comme suit , pour former **ABCD** j'ai besoin de **ABC** , **ABD** , **ACD** et **BCD** en d'autres termes, nous pourrions dire $ABC \cup ABD \cup \dots \cup BCD = ABCD$.

Pour former **AB** j'ai besoin de **ABC** et **ABD** .

Pour former **A** j'ai besoin de **AB** , **AC** et **AD** .

Pour former **all** j'ai besoin de **A** , **B** , **C** et **D** .

Si on part de la feuille le graphe se lit **all** est constitué de **A** , **B** , **C** et **D** etc en d'autres termes, nous pourrions dire que $all \cap A = A$ (ou *all GROUP BY A* donne **A**) .

Dans l'approche **BUC** autrement dit "Bottom up Computation" c'est la dernière formulation que nous allons utiliser .

2.3 Exemple :

Reprenons l'exemple précédemment évoqué , avec **A** : le nom du produit , **B** : la saison et **C** : la localisation . Pour chacune de ces 3 dimensions nous nous restreindrons aux valeurs les plus fréquentes par exemple : *écran-télé-type-X* , *écran-télé-type-Y* et *écran-ordinateur-type-X* , *saison-hiver* , *saison-été* et *saison-printemps* , *USA* , *Europe* et *Asie* .

Nous constatons que les saisons et les localisations sont toutes présentes en effet le *threshold* n'a pas exclu ces valeurs car elles sont toutes significatives pour l'analyse du marché en revanche les produits ne sont pas tous tendances (significatives) .

2.4 Implémentation :

Le **BUC** comme nous l'avons dit commence par l'entité la plus agrégée vers l'entité la moins agrégée. (du plus grand vers le plus petit/ du moins de *GROUP BY* vers le plus de *GRP BY*)

Procedure BottomUpCube(input, dim)

Inputs:
input: The relation to aggregate.
dim: The starting dimension for this iteration.

Globals:
constant numDims: The total number of dimensions.
constant cardinality[numDims]: The cardinality of each dimension.
constant minsup: The minimum number of tuples in a partition for it to be output.
outputRec: The current output record.
dataCount[numDims]: Stores the size of each partition.
dataCount[i] is a list of integers of size cardinality[i].

Outputs:
One record that is the aggregation of input.
Recursively, outputs CUBE(dim, ..., numDims) on input (with minimum support).

Method:
1: Aggregate(input); // Places result in outputRec
2: **if** input.count() == 1 **then** // Optimization
 WriteAncestors(input[0], dim); **return**;
3: write outputRec;
4: **for** d = dim ; d < numDims ; d++ **do**
5: let C = cardinality[d];
6: Partition(input, d, C, dataCount[d]);
7: let k = 0;
8: **for** i = 0 ; i < C ; i++ **do** // For each partition
9: let c = dataCount[d][i]
10: **if** c >= minsup **then** // The BUC stops here
11: outputRec.dim[d] = input[k].dim[d];
12: BottomUpCube(input[k ... k+c], d+1);
13: **end if**
14: k += c;
15: **end for**
16: outputRec.dim[d] = ALL;
17: **end for**

Figure 5: Algorithm BottomUpCube (BUC)

2.4.1 Paramètres :

Paramètre d'entrée :

Dans le pseudo code ci-dessus nous avons deux paramètres d'entrés :

Input:

L'input peut pour l'instant être vue comme un paramètre composite : un dictionnaire qui associe des clés les identifiants des colonnes aux valeurs de chaque tuple pour chacune de ces colonnes .

Ainsi nous avons un input avec par exemple une ligne/tuple :

$$[A : val_a1 , B : val_b2 , C : val_c2 , D : val_d1]$$

Dim :

Dim quant à lui décrit l'index de la dimension actuellement processor , en d'autres termes la dimension pour laquelle on parcourt l'ensemble de ses valeurs .

Paramètre de sorties :

Le résultat de l'opération *CUBE* affiché petit à petit

2.4.2 Commentaire sur l'Algorithme :

Premièrement, on agrège l'input (à la manière d'un *GROUP BY*) on confond les entrées/lignes qui partagent la même valeur sur la dimension concernée (*index_{dim}*) .

Ensuite on parcourt l'ensemble des valeurs distinctes sur laquelle on a agrégé la dimension d'indexe *index_{dim}* pour chacune de ces valeurs on partitionne les lignes sur la dimension suivante .

On parcourt l'ensemble de ces partitions (i de 0 à n) pour chacune on teste si sa cardinalité (le nombre de ligne) est supérieure à la valeur *threshold* si oui on rappelle **BUC** sur *partition[i]* et la dimension suivante : *k+1* .

Ce processus est itéré pour chacune des permutations du mot **ABCD** (représentant les 4 dimensions de notre exemple) .

2.4.3 Nota Bene :

Les mots **ABCD** , **ACDB** , **BCDA** ect sont équivalents pour preuve :

$$GRP\ BY\ A,B,C,D \leftrightarrow GRP\ BY\ B,C,D,A \text{ (cf. l'opération } CUBE \text{)}.$$

Ainsi nous avons : $\frac{n(n+1)}{2}$ computation possibles .

2.5 Implémentation de l'Algorithme :

BUC-implémentation <https://github.com/LonyI175/BUC-Iceberg-Cube> forked from AsaadMe repository .

Pour une première implémentation de l'algorithme il est pratique de le pensé comme une chaine de chiffre binaire 0,1,* (avec * représentant le caractère joker) .

Tout d'abord, nous devons binariser les entrées pour ce fait nous partagerons chacune des dimensions en 2 parties :

- $A[i] \leq sup_a \implies 1$
- $A[i] > sup_a \implies 0$

Une fois ceci fait pour chacune des dimensions :

```
01 | def pre_process(input_file_name: str) -> list[dict]:
02 |     """Convert records to dict: {"A": 0 or 1, ...}"""
03 |
04 |     out_list = []
05 |
06 |     with open(input_file_name, "r") as dataset:
07 |
08 |         for line in dataset:
09 |             out_line = {}
10 |
11 |             if re.search(r"\b[0-3]\b", line):
12 |
13 |                 out_line["A"] = 1
14 |             else:
15 |                 out_line["A"] = 0
```

Nous allons ensuite merger ces partitions sur le nouvel indice de dimension (nous regroupons les entrées qui ont la même valeur sur la dimension identifiée par l'indice : *index_dim*) .

```
01 | for row in input:
02 |     if row.get(cur_dim) == 1:
03 |         aggr_list_1.append(row)
04 |     else:
05 |         aggr_list_0.append(row)
```

Ensuite, nous calculons le nombre de ligne dans chacune de ces partitions si ce nombre est supérieur aux *threshold* alors cette partition est significative et elle fait l'office d'un nouvel appel de l'algorithme **BUC** sinon elle est discarded :

```

01 | datacount = [len(aggr_list_0), len(aggr_list_1)]
02 |
03 |     if datacount[1] >= min_sup:
04 |
05 |         cell_copy = cell.copy()
06 |         cell_copy[cur_dim] = "1"
07 |         print("counter : "+str(my_counter)+ f" {cell_repr(cell_copy)}: {
datacount[1]}")
08 |
09 |         if (index:= dims.index(cur_dim) + 1) < len(dims):
10 |             BUC(aggr_list_1, index, cell_copy, my_counter)

```

Note :

Seul le traitement de la liste *aggr_list_1* est présenté , mais il en va de même pour la liste concernant le second intervalle : *aggr_list_0* .

2.5.1 Résultat :

Fichier utilisé :

Database <https://github.com/LonyI175/BUC-Iceberg-Cube/blob/master/Dataset.txt> (réalisé par Mehran Asaad)

Aperçu :

```

cur_dim:A counter : 1 parent :0
[{'A': 0, 'B': 1, 'C': 0, 'D': 0, 'E': 0}, {'A': 0, 'B': 1, 'C': 1, 'D': 0, 'E': 1}, {'A': 0, 'B': 1, 'C': 1, 'D': 1,
0 : (*, *, *, *, *): 13610
1 : (*, *, *, *, *): 86390
counter : 1 (1, *, *, *, *): 86390
cur_dim:B counter : 2 parent :1
[{'A': 1, 'B': 0, 'C': 1, 'D': 1, 'E': 1}, {'A': 1, 'B': 0, 'C': 1, 'D': 1, 'E': 1}, {'A': 1, 'B': 0, 'C': 0, 'D': 1,
0 : (1, *, *, *, *): 15328
1 : (1, *, *, *, *): 71062
counter : 2 (1, 1, *, *, *): 71062
cur_dim:C counter : 3 parent :2
[{'A': 1, 'B': 1, 'C': 0, 'D': 1, 'E': 1}, {'A': 1, 'B': 1, 'C': 0, 'D': 1, 'E': 1}, {'A': 1, 'B': 1, 'C': 0, 'D': 1,
0 : (1, 1, *, *, *): 8716
1 : (1, 1, *, *, *): 62346
counter : 3 (1, 1, 1, *, *): 62346
cur_dim:D counter : 4 parent :3
[{'A': 1, 'B': 1, 'C': 1, 'D': 0, 'E': 1}, {'A': 1, 'B': 1, 'C': 1, 'D': 0, 'E': 0}, {'A': 1, 'B': 1, 'C': 1, 'D': 0,
0 : (1, 1, 1, *, *): 7129
1 : (1, 1, 1, *, *): 55217
counter : 4 (1, 1, 1, 1, *): 55217
cur_dim:E counter : 5 parent :4
[{'A': 1, 'B': 1, 'C': 1, 'D': 1, 'E': 0}, {'A': 1, 'B': 1, 'C': 1, 'D': 1, 'E': 0}, {'A': 1, 'B': 1, 'C': 1, 'D': 1,
0 : (1, 1, 1, 1, *): 6428
1 : (1, 1, 1, 1, *): 48789
counter : 5 (1, 1, 1, 1, 1): 48789
counter : 5 (1, 1, 1, 1, 0): 6428
counter : 4 (1, 1, 1, 0, *): 7129
cur_dim:E counter : 6 parent :4
[{'A': 1, 'B': 1, 'C': 1, 'D': 0, 'E': 0}, {'A': 1, 'B': 1, 'C': 1, 'D': 0, 'E': 0}, {'A': 1, 'B': 1, 'C': 1, 'D': 0,
0 : (1, 1, 1, 0, *): 1039
1 : (1, 1, 1, 0, *): 6090

```

References

- [1] *BUC implémentation* , Mehran Asaad , 2021, <https://github.com/LonyI175/BUC-Iceberg-Cube>
- [2] *Bottom-Up Computation of Sparse and Iceberg CUBES* , Kevin Beyer and Raghu Ramakrishnan , 1999, <https://dl.acm.org/doi/10.1145/304181.304214>
- [3] *Iceberg-Cubes* , Uregina , 2005, <http://www2.cs.uregina.ca/~dbd/cs831/notes/dcubes/iceberg.html>