# SQL Aggregation Operations for Horizontal Distributed Data

## Mouni Kurama[#1], K. Johnpaul[#2]

#1 CSE, Nova College of Engineering & Technology, Vegavaram, Jangareddy Gudem,
#2B-Tech, M-Tech, Associate Professor, Nova College of Engineering & Technology, Vegavaram, Jangareddy Gudem

**Abstract:** *Importance of information sets could be ascribed in distinctive spaces, for example, Verification of productions, longitudinal exploration, Interdisciplinary utilization of information, Valorization and so on. As opposed to utilizing outsider business devices to create information sets for exploration from RDBMS, SQL Built In Aggregates might be utilized. Fundamental SQL collections restrictions to give back one section for every collected gathering utilizing gathering capacities is overcome by a straightforward, yet influential, methods(case,pivot,spj) to create amassed sections in a flat even format helped with an agreeable programming dialect. CASE and PIVOT assessment routines are essentially speedier than the SPJ strategy, considering the imperativeness of an all encompassing execution to improve the working of existing local RDBMS systems, for example, SPJ as opposed to manufacture new ones, we propose to utilize Join Enumeration methodologies to upgrade the execution of SPJ. The methods incorporates a question tree era with quantifiers calculation, which incorporates relations referenced by the join predicate that are utilized to partner each one join predicate furthermore considering extra relations required by a predicate to save the semantics of the first inquiry. The methods enhances SPJ execution fundamentally since applying totals ahead of schedule in inquiry preparing can give huge execution upgrades.*

**Index Terms:** *Horizontal Aggregation, CASE, PIVOT,SQL Data.*

## 1. INTRODUCTION

There are two primary parts in such SQL code: joins and collections. The most generally known accumulation is the total of a section over gatherings of lines. There exist numerous conglomeration capacities and administrators in SQL. Tragically, all these accumulations have confinements to manufacture information sets for information mining purposes. The primary reason is that, as a rule, information sets that are put away in a social database (or an information distribution center) originate from On-Line Transaction Processing (OLTP) frameworks where database blueprints are exceptionally standardized. In view of current accessible capacities and statements in SQL, a noteworthy exertion is obliged to process collections. Such exertion is because of the sum and multifaceted nature of SQL code that needs to be composed, advanced and tried. Standard conglomerations are difficult to decipher when there are numerous result lines. new class of total capacities that total numeric declarations and transpose results to create an information set with a level format. Capacities having a place with this class are called flat aggregations. first, they speak to a layout to produce SQL code from an information mining apparatus. This SQL code decreases manual work in the information planning stage in an information mining undertaking. Second, since SQL code is naturally produced it is liable to be more productive than SQL code composed by an end client. Third, the information set might be made altogether inside the DBMS. Even accumulations simply oblige a little punctuation expansion to total capacities brought in a SELECT explanation.

We create a procedure for pushing Gps down inquiry trees of Select-task join may utilize conglomerations like max, whole, and so on and that utilize discretionary capacities as a part of their choice conditions. Our procedure pushes down to the most reduced levels of a question tree total calculation, copy end, and capacity reckoning.

## 2. HORIZONTAL AGGREGATIONS

Our main goal is to define a template to generate SQL code combining aggregation and transposition (pivoting). A second goal is to extend the SELECT statement with a clause that combines transposition with aggregation. A method, SPJ method, is used to evaluate horizontal aggregations which relies on relational operations. That is, select, project, join and aggregation queries. In order to

evaluate this query the query optimizer takes three input parameters: (1) the input table F, (2) the list of grouping columns $L_1;\ldots;L_m$, (3) the column to aggregate (A). In a horizontal aggregation there are four input parameters to generate SQL code: 1) the input table F, 2) the list of GROUP BY columns $L_1;L_j$, 3) the column to aggregate (A), 4) the list of transposing columns $R_1;\ldots;R_k$.

SELECT $L_1;\ldots;L_J$, H(A BY $R_1;\ldots;R_k$)

FROM F

GROUP BY $L_1;\ldots;L_J$;

The result rows are determined by columns $L_1;\ldots;L_J$ in the GROUP BY clause if present. Result columns are determined by all potential combinations of columns $R_1;\ldots;R_k$, where k = 1 is the default.

The main reasons are that any insertion into F during evaluation may cause inconsistencies: (1) it can create extra columns in $F_H$, for a new combination of $R_1;\ldots;R_k$; (2) it may change the number of rows of $F_H$, for a new combination of $L_1;\ldots;L_J$; (3) it may change actual aggregation values in $F_H$.

Therefore, the result table FH must have as primary key the set of grouping columns $\{L_1;\ldots;L_J\}$ and as non-key columns all existing combinations of values $R_1;\ldots;R_k$.

A horizontal aggregation exhibits the following properties:

1) n= $|F_H|$ matches the number of rows in a vertical aggregation grouped by $L1;\ldots;Lj$.

2) d = $|\pi_{R1,\ldots,Rk}(F)|$

3) Table $F_H$ may potentially store more aggregated values than $F_V$ due to nulls. That is, $|F_V|\leq$ nd.

*DBMS limitations:* On the other hand, the second important issue is automatically generating unique column names. However, these are not important limitations because if there are many dimensions that is likely to correspond to a sparse matrix (having many zeroes or nulls) on which it will be difficult or impossible to compute a data mining model. The column name length issue can be solved by generating column identifiers with integers and creating a description table that maps identifiers to full descriptions, but the meaning of each dimension is lost. An alternative is the use of abbreviations, which may require manual input.

## 3. SPJ METHOD

The basic idea is to create one table with a vertical aggregation for each result column, and then join all those tables to produce FH. We aggregate from F into d projected tables with d Select-Project-Join-Aggregation queries (selection, projection, join, aggregation). Each table $F_I$ corresponds to one subgrouping combination and has $\{L1;\ldots;Lj\}$ as primary key and an aggregation on A as the only non-key column. It is necessary to introduce an additional table $F_0$, that will be outer joined with projected tables to get a complete result set. We propose two basic sub-strategies to compute $F_H$. The first one directly aggregates from F. The second one computes the equivalent vertical aggregation in a temporary table $F_V$ grouping by $L1;\ldots;Lj;R_1;\ldots;R_k$.



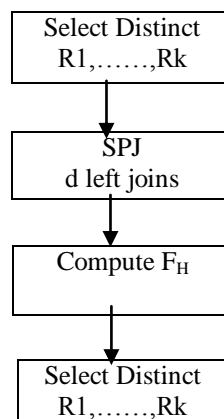**Figure 1.** *Main steps of methods based on FV (un optimized)*

The statement to compute FV gets a cube:

INSERT INTO $F_V$

SELECT $L_1; \ldots ; L_J; R_1; \ldots ; R_k \, V(A)$

FROM F

GROUP BY $L_1; \ldots ; L_J; R_1; \ldots ; R_k;$

Table F0 de_nes the number of result rows, and builds the primary key. F0 is populated so that it contains every existing combination of $L_1; \ldots ; L_J$. Table F0 has { $L_1; \ldots ; L_J$ } as primary key and it does not have any non-key column.

INSERT INTO $F_0$

SELECT DISTINCT $L_1; \ldots ; L_J$

FROM $\{F | F_V \}$;

In the following discussion I $\in$ {1;… ; d}. we use h to make writing clear, mainly to define boolean expressions. We need to get all distinct combinations of subgrouping columns $R_1; \ldots ; R_k$, to create the name of dimension columns, to get d, the number of dimensions, and to generate the boolean expressions for WHERE clauses. Each WHERE clause consists of a conjunction of k equalities based on $R_1 ; \ldots ; Rk$.

SELECT DISTINCT $R_1; \ldots ; R_k$

FROM $\{F|F_V\}$;

Tables $F_1; \ldots ; F_d$ contain individual aggregations for each combination of $R_1; \ldots ; R_k$. The primary key of table FI is { $L_1; \ldots ; L_J$ }.

INSERT INTO $F_I$

SELECT $L1; \ldots ; Lj ; V(A)$

FROM {F|FV}

WHERE R1 = v1I AND .. AND Rk = vkI

GROUP BY $L1; \ldots ; Lj$ ;

Then each table $F_I$ aggregates only those rows that correspond to the Ith unique combination of $R_1; \ldots ; R_k$, given by the WHERE clause. A possible optimization is synchronizing table scans to compute the d tables in one pass. Finally, to get $F_H$ we need d left outer joins with the d + 1 tables so that all individual aggregations are properly assembled as a set of d dimensions for each group. Outer joins set result columns to null for missing combinations for the given group. In general, nulls should be the default value for groups with missing combinations. We believe it would be incorrect to set the result to zero or some other number by default if there are no qualifying rows. Such approach should be considered on a per-case basis.

INSERT INTO $F_H$
SELECT
$F_0.L_1; F_0.L_2; \ldots ; F_0.L_J;$
$F_1.A; F_2.A; \ldots ; F_d.A$
FROM $F_d$
LEFT OUTER JOIN $F_1$
ON $F_0.L_1 = F_1.L_1$ and … and $F_0.L_J = F_1.L_J$
LEFT OUTER JOIN $F_2$
ON $F_0.L_1 = F_2.L_1$ and … and $F_0.L_J = F_2.L_J$
….
LEFT OUTER JOIN $F_d$
ON $F_0.L_1 = F_d.L_1$ and …. and $F_0.L_J = F_d.L_J$ ;

We introduce the notion of a generalized projection that unifies duplicate eliminating projections corresponds to the SQL distinct adjective, duplicate preserving projections, group by, and aggregations, in a common framework.

## 4. GENERALIZED PROJECTION

We introduce a generalized projection operator, denoted by the symbol $\pi$ , that is similar to aggregation operator. A GP takes as its argument a relation R and outputs a new relation based on the subscript of the GP. The subscript specifies the computation to be done on R. The subscript has two parts:

1. A set of group by components. We refer to them as components and not attribute because they may be functions of attributes and not just attributes. For instance, the GP $\pi_{A*B}$ (R) is written as the following SQL query:

Select (A*B) from R group by (A*B).

2. A set of aggregate components. For example, we can write the GP $\pi_{D,max(S)}$ (R) as the query:

Select D, max(S) from R group by D.

Here D is the only group by component and max(S) is the only aggregate component. It is simple to observe that a GP has exactly one tuple for each value of the group by components and thus does not produce any duplicates in its output. Here class of queries expressed in a query tree. The permitted query trees have ve types of nodes: selection nodes, projection nodes, cross-product nodes, group by nodes, and aggregate-group by node pairs.

Projections may preserve duplicates or discard them.. Selection nodes eliminate tuples from the input relation, group by nodes do projection duplicate elimination, and cross-product nodes output the cross product of two input relations. Aggregate- group by node pairs have a group by node followed by an aggregate node. An aggregate-group by node pair produces as output a relation with one tuple for every distinct value in the input relation of the group by attributes.

GPs are incorporated into query trees using a two step process:

1. Push GPs down a query tree and annotate the query tree with a GP above each node in the tree.

2. Rewrite the annotated query tree to incorporate the GPs that the query optimizer chooses to evaluate and to eliminate all other GPs introduced in the push-down process.

After the top-down pass associates a GP with some or all nodes of the query tree, the query optimizer decides which GPs improve the query plan. The other GPs are removed from the tree.

## 5. PERFORMANCE RESULTS

Most inquiries are not intrigued by individual tuples of this connection, yet rather total properties of this connection. Consequently as a rule, we have to do a groupby on a non-key property of this connection. At the point when this connection is joined with some other connection,  that need not be collected.

```
Algorithm 1 the construction of ES

GetES
Input: original query tree G = (V,E) with quantifiers q={q₁,…,
q_N }
Output:  a set contains the NS and ES  of each join predicate
1   Loop          //for each join predicate p
2     if ○p is inner join or antijoin
3       for each relation r ref(p)
4          s=s+ set_outerjoin(r);
5       for each member of s
6          set the set_outerjoin of the member to be s
7     if ○p is outerjoin
8       for each relation r in ref(nullproducting(p))
9          v=v+set_outerjoin(r);
10        set the ES of p to be v
11        for each relation r in ref(preserving(p))
12           u=u+ set_antijoin(r);
13        for each relation r ref(nullproducting(p))
14           add all the tables in u to the set_antijoin(r);
15    if ○p is antijoin
16        for each r in ref(preserving(p))
17           w=w+ set_antijoin(r);
18        set the ES of p to be w
```

**Figure 2.** *SPJ Performance evaluation algorithm*

In such cases, our system would diminish significantly the span of the enormous table before we did a join. It could be contended that in such cases a join calculation like a hash join could be utilized to attain a comparable result. In any case, hash joins are dificult to execute in practice and not generally actualized. Single table accumulations being an ordinarily utilized gimmick of SQL exist within generally frameworks. Our streamlining, when connected to question plans, possibly meddles with join requesting, since we lessen the span of the relations partaking in the join.

## 6. CONCLUSION

Level accumulations help in choice making giving a naturally visible perspective of whole business. Utilizes a basic, yet effective, mining routines (CASE, PIVOT, SPJ) of RDBMS to create accumulated segments in an even plain design. Execution of CASE, PIVOT was viewed as average regarding rate and scalability.spj slacks in velocity and versatility measurements. It is vital to upgrade the working of existing local RDBMS strategies, for example, SPJ as opposed to construct new ones to fuse mining. So we propose Join Enumeration methods. The methodologies incorporates a question tree era with quantifiers calculation, which incorporates relations referenced by the join predicate that are utilized to partner each one join predicate furthermore considering extra relations required by a predicate to protect the semantics of the first inquiry.

## REFERENCES

[1] A. Witkowski, S. Bellamkonda, T. Bozkaya, G. Dorman, N. Folkert, A. Gupta, L. Sheng, and S. Subramanian. Spreadsheets in RDBMS for OLAP. In *Proc. ACM SIGMOD Coference*, pages 52.63, 2003.

[2] Venky Harinarayan ,Ashish Guptay "Generalized Projections: a Powerful Query-Optimization Technique "

[3] "Vertical and Horizontal Percentage Aggregations", Carlos Ordonez Teradata, NCR San Diego, CA 92127, USA.

[4] G. Bhargava, P. Goel, and B.R. Iyer. Hypergraph based reordering of outer join queries with complex predicates. In *ACM SIGMOD Conference*, pages 304.315, 1995.

[5] U. Dayal, N. Goodman, and R. H. Katz. "An Extended Relational Algebra with Control over Duplicate Elimination". In Proceedings of the ACM Symposium on Principles of Database Systems, 1982, pages 117-123.

[6] Venky Harinarayan and Ashish Gupta. Optimization Using Tuple Subsumption. To appear in ICDT 95, January 1995.