

Advanced Aggregate Functions & Advanced Grouping Operations

Data Transformation and Aggregation: Solutions

May 19, 2025

Contents

1	Category: Advanced Aggregate Functions	6
1.1	STRING_AGG(expression, separator [ORDER BY ...])	6
1.1.1	Practice Meaning, Values, Relations, Advantages	6
1.1.2	Practice Disadvantages	6
1.1.3	Practice Inefficient Alternatives Avoidance	6
1.2	ARRAY_AGG(expression [ORDER BY ...])	7
1.2.1	Practice Meaning, Values, Relations, Advantages	7
1.2.2	Practice Disadvantages	7
1.2.3	Practice Inefficient Alternatives Avoidance	8
1.3	JSON_AGG(expression [ORDER BY ...])	8
1.3.1	Practice Meaning, Values, Relations, Advantages	8
1.3.2	Practice Disadvantages	8
1.3.3	Practice Inefficient Alternatives Avoidance	9
1.4	PERCENTILE_CONT(fraction) WITHIN GROUP (ORDER BY sort_expression)	9
1.4.1	Practice Meaning, Values, Relations, Advantages	9
1.4.2	Practice Disadvantages	9
1.4.3	Practice Inefficient Alternatives Avoidance	10
1.5	CORR(Y, X)	10
1.5.1	Practice Meaning, Values, Relations, Advantages	10
1.5.2	Practice Disadvantages	10
1.5.3	Practice Inefficient Alternatives Avoidance	11
1.6	REGR_SLOPE(Y, X)	11
1.6.1	Practice Meaning, Values, Relations, Advantages	11
1.6.2	Practice Disadvantages	11
1.6.3	Practice Inefficient Alternatives Avoidance	12
2	Category: Advanced Grouping Operations	13
2.1	GROUPING SETS ((set1), (set2), ...)	13
2.1.1	Practice Meaning, Values, Relations, Advantages	13
2.1.2	Practice Disadvantages	13
2.1.3	Practice Inefficient Alternatives Avoidance	13
2.2	ROLLUP (col1, col2, ...)	14
2.2.1	Practice Meaning, Values, Relations, Advantages	14

2.2.2	Practice Disadvantages	14
2.2.3	Practice Inefficient Alternatives Avoidance	15
2.3	CUBE (col1, col2, ...)	15
2.3.1	Practice Meaning, Values, Relations, Advantages	15
2.3.2	Practice Disadvantages	15
2.3.3	Practice Inefficient Alternatives Avoidance	16
3	Hardcore Combined Problem	17

Global Dataset for PostgreSQL

The following SQL code creates and populates the necessary tables for the exercises. This is the same dataset provided with the exercises.

```
1  -- Drop tables if they exist to ensure a clean setup
2  DROP TABLE IF EXISTS EmployeeProjects CASCADE;
3  DROP TABLE IF EXISTS Sales CASCADE;
4  DROP TABLE IF EXISTS Employees CASCADE;
5  DROP TABLE IF EXISTS Departments CASCADE;
6  DROP TABLE IF EXISTS Projects CASCADE;
7  DROP TABLE IF EXISTS Products CASCADE;
8  DROP TABLE IF EXISTS Regions CASCADE;
9
10 -- Table: Departments
11 CREATE TABLE Departments (
12     departmentId INT PRIMARY KEY,
13     departmentName VARCHAR(100) NOT NULL,
14     locationCity VARCHAR(50)
15 );
16
17 -- Table: Employees
18 CREATE TABLE Employees (
19     employeeId INT PRIMARY KEY,
20     firstName VARCHAR(50) NOT NULL,
21     lastName VARCHAR(50) NOT NULL,
22     email VARCHAR(100) UNIQUE,
23     hireDate DATE NOT NULL,
24     salary DECIMAL(10, 2) NOT NULL,
25     departmentId INT,
26     managerId INT,
27     performanceScore NUMERIC(3,2), -- Score from 0.00 to 5.00
28     skills TEXT[], -- Array of skills
29     FOREIGN KEY (departmentId) REFERENCES Departments(departmentId),
30     FOREIGN KEY (managerId) REFERENCES Employees(employeeId)
31 );
32
33 -- Table: Projects
34 CREATE TABLE Projects (
35     projectId INT PRIMARY KEY,
36     projectName VARCHAR(100) NOT NULL,
37     startDate DATE,
38     deadlineDate DATE,
39     budget DECIMAL(12,2)
40 );
41
42 -- Table: EmployeeProjects
43 CREATE TABLE EmployeeProjects (
44     assignmentId SERIAL PRIMARY KEY,
45     employeeId INT,
46     projectId INT,
47     hoursWorked INT,
48     taskNotes TEXT,
49     FOREIGN KEY (employeeId) REFERENCES Employees(employeeId),
50     FOREIGN KEY (projectId) REFERENCES Projects(projectId)
51 );
52
53 -- Table: Regions
54 CREATE TABLE Regions (
55     regionId INT PRIMARY KEY,
56     regionName VARCHAR(50) NOT NULL UNIQUE
57 );
58
59 -- Table: Products
60 CREATE TABLE Products (
61     productId INT PRIMARY KEY,
62     productName VARCHAR(100) NOT NULL,
63     category VARCHAR(50),
64     standardCost DECIMAL(10, 2),
65     listPrice DECIMAL(10, 2)
66 );
67
```

```

68 -- Table: Sales
69 CREATE TABLE Sales (
70     saleId INT PRIMARY KEY,
71     productId INT,
72     employeeId INT,
73     saleDate DATE NOT NULL,
74     quantity INT NOT NULL,
75     regionId INT,
76     notes JSONB, -- e.g., {"customerSatisfaction": 5, "followUpRequired": true}
77     FOREIGN KEY (productId) REFERENCES Products(productId),
78     FOREIGN KEY (employeeId) REFERENCES Employees(employeeId),
79     FOREIGN KEY (regionId) REFERENCES Regions(regionId)
80 );
81
82 -- Insert data into Departments
83 INSERT INTO Departments (departmentId, departmentName, locationCity) VALUES
84 (1, 'Human Resources', 'New York'),
85 (2, 'Engineering', 'San Francisco'),
86 (3, 'Sales', 'Chicago'),
87 (4, 'Marketing', 'New York'),
88 (5, 'Research', 'San Francisco');
89
90 -- Insert data into Employees
91 INSERT INTO Employees (employeeId, firstName, lastName, email, hireDate, salary,
92     departmentId, managerId, performanceScore, skills) VALUES
93 (101, 'Alice', 'Smith', 'alice.smith@example.com', '2020-01-15', 70000, 2, NULL, 4.50,
94     ARRAY['Java', 'Python', 'SQL']),
95 (102, 'Bob', 'Johnson', 'bob.johnson@example.com', '2019-03-01', 80000, 2, 101, 4.20,
96     ARRAY['Python', 'Machine Learning']),
97 (103, 'Carol', 'Williams', 'carol.williams@example.com', '2021-07-30', 60000, 1, NULL,
98     3.90, ARRAY['HR Policies', 'Recruitment']),
99 (104, 'David', 'Brown', 'david.brown@example.com', '2018-06-11', 95000, 2, 101, 4.80,
100     ARRAY['Java', 'Spring Boot', 'Microservices']),
101 (105, 'Eve', 'Davis', 'eve.davis@example.com', '2022-01-10', 75000, 3, NULL, 4.10, ARRAY
102     ['Salesforce', 'Negotiation']),
103 (106, 'Frank', 'Miller', 'frank.miller@example.com', '2019-11-05', 120000, 3, 105, 4.60,
104     ARRAY['Key Account Management', 'CRM']),
105 (107, 'Grace', 'Wilson', 'grace.wilson@example.com', '2020-08-20', 65000, 4, NULL, 3.70,
106     ARRAY['SEO', 'Content Creation']),
107 (108, 'Henry', 'Moore', 'henry.moore@example.com', '2023-02-18', 55000, 1, 103, 4.00,
108     ARRAY['Onboarding', 'Employee Relations']),
109 (109, 'Ivy', 'Taylor', 'ivy.taylor@example.com', '2017-05-25', 110000, 5, NULL, 4.90,
110     ARRAY['Research Methodologies', 'Statistical Analysis', 'Python']),
111 (110, 'Jack', 'Anderson', 'jack.anderson@example.com', '2021-10-01', 72000, 5, 109,
112     4.30, ARRAY['Lab Techniques', 'Data Analysis']),
113 (111, 'Kevin', 'Spacey', 'kevin.spacey@example.com', '2020-05-15', 65000, 4, 107, 4.1,
114     ARRAY['Digital Marketing', 'Analytics']),
115 (112, 'Laura', 'Palmer', 'laura.palmer@example.com', '2021-08-01', 90000, 5, 109, 4.7,
116     ARRAY['Quantum Physics', 'Research']),
117 (113, 'Dale', 'Cooper', 'dale.cooper@example.com', '2019-09-10', 130000, 3, 105, 4.8,
118     ARRAY['Strategic Sales', 'Leadership']),
119 (114, 'Audrey', 'Horne', 'audrey.horne@example.com', '2022-03-20', 60000, 1, 103, NULL,
120     ARRAY['Payroll', 'Conflict Resolution']);
121
122 -- Insert data into Projects
123 INSERT INTO Projects (projectId, projectName, startDate, deadlineDate, budget) VALUES
124 (1, 'Alpha Platform', '2023-01-01', '2023-12-31', 500000),
125 (2, 'Beta Feature', '2023-03-15', '2023-09-30', 150000),
126 (3, 'Gamma Initiative', '2023-06-01', '2024-05-31', 750000),
127 (4, 'Delta Rollout', '2022-11-01', '2023-07-30', 300000);
128
129 -- Insert data into EmployeeProjects
130 INSERT INTO EmployeeProjects (employeeId, projectId, hoursWorked, taskNotes) VALUES
131 (101, 1, 120, 'Developed core APIs'),
132 (102, 1, 100, 'Machine learning model integration'),
133 (104, 1, 150, 'Backend services for Alpha'),
134 (101, 2, 80, 'API refinement for Beta feature'),
135 (105, 3, 200, 'Sales strategy for Gamma'),
136 (106, 3, 180, 'Client acquisition for Gamma'),
137 (107, 4, 90, 'Marketing campaign for Delta'),
138 (109, 2, 110, 'Research for Beta feature improvements'),
139 (110, 2, 70, 'Data analysis for Beta feature testing'),

```

```

126 (102, 3, 50, 'Consulting on ML aspects for Gamma');
127
128 -- Insert data into Regions
129 INSERT INTO Regions (regionId, regionName) VALUES
130 (1, 'North'), (2, 'South'), (3, 'East'), (4, 'West'), (5, 'Central');
131
132 -- Insert data into Products
133 INSERT INTO Products (productId, productName, category, standardCost, listPrice) VALUES
134 (1, 'Laptop Pro', 'Electronics', 800, 1200),
135 (2, 'Smartphone X', 'Electronics', 400, 700),
136 (3, 'Office Chair', 'Furniture', 100, 250),
137 (4, 'Desk Lamp', 'Furniture', 20, 45),
138 (5, 'Software Suite', 'Software', 50, 150),
139 (6, 'Advanced CPU', 'Components', 250, 400),
140 (7, 'Graphics Card', 'Components', 300, 550);
141
142 -- Insert data into Sales
143 INSERT INTO Sales (saleId, productId, employeeId, saleDate, quantity, regionId, notes)
144 VALUES
145 (1, 1, 105, '2022-01-20', 2, 1, '{"customerSatisfaction": 5, "followUpRequired": false}'
146 ),
147 (2, 2, 106, '2022-02-10', 5, 2, '{"customerSatisfaction": 4, "discountApplied": "10%"}'),
148 (3, 1, 105, '2022-02-15', 1, 1, '{"customerSatisfaction": 4, "followUpRequired": true, "
149 feedback": "Needs faster shipping options"}'),
150 (4, 3, 106, '2022-03-05', 10, 3, NULL),
151 (5, 4, 105, '2023-03-22', 20, 4, '{"customerSatisfaction": 3}'),
152 (6, 5, 106, '2023-04-10', 50, 1, '{"customerSatisfaction": 5, "bulkOrder": true}'),
153 (7, 2, 105, '2023-04-18', 3, 2, '{"customerSatisfaction": 5}'),
154 (8, 1, 106, '2022-05-01', 2, 3, '{"notes": "Repeat customer"}'),
155 (9, 3, 105, '2022-05-25', 8, 4, NULL),
156 (10, 5, 106, '2023-06-11', 30, 5, '{"customerSatisfaction": 4, "followUpRequired": true}
157 '),
158 (11, 6, 102, '2023-07-01', 5, 1, '{"source": "Tech Expo"}'),
159 (12, 7, 104, '2023-07-05', 3, 2, '{"source": "Internal Purchase"}'),
160 (13, 1, 105, '2022-01-25', 3, 1, '{"customerSatisfaction": 5}'),
161 (14, 2, 105, '2023-02-12', 2, 2, '{"customerSatisfaction": 3, "feedback": "Item was
162 backordered"}'),
163 (15, 1, 106, '2023-01-30', 1, 1, NULL),
164 (16, 3, 113, '2022-08-15', 12, 2, '{"customerSatisfaction": 5}'),
165 (17, 4, 105, '2022-09-01', 25, 3, '{"customerSatisfaction": 4, "notes": "Urgent delivery
166 "}''),
167 (18, 5, 106, '2023-08-20', 60, 4, '{"bulkOrder": true}'),
168 (19, 6, 113, '2023-09-05', 8, 5, NULL),
169 (20, 7, 105, '2023-10-10', 4, 1, '{"customerSatisfaction": 5, "followUpRequired": true}'
170 );
171
172 -- Update data for NULL examples
173 UPDATE Employees SET departmentId = NULL WHERE employeeId = 108; -- Henry Moore has no
174 department
175 UPDATE Sales SET regionId = NULL WHERE saleId = 4; -- Sale 4 has no region
176 UPDATE Products SET category = NULL WHERE productId = 4; -- Desk Lamp has no category

```

Listing 1: Global Dataset for Exercises (Identical to Exercises Document)

1 Category: Advanced Aggregate Functions

1.1 STRING_AGG(expression, separator [ORDER BY ...])

1.1.1 Practice Meaning, Values, Relations, Advantages

```
1 SELECT
2     COALESCE(d.departmentName, 'No Department Assigned') AS
    departmentName,
3     STRING_AGG(e.firstName, ', ' ORDER BY e.firstName) AS employeeNames
4 FROM Employees e
5 LEFT JOIN Departments d ON e.departmentId = d.departmentId
6 GROUP BY d.departmentId, d.departmentName
7 ORDER BY d.departmentName NULLS FIRST;
```

Listing 2: Solution 1.1.1: Department employee list

1.1.2 Practice Disadvantages

```
1 -- Disadvantage: STRING_AGG can produce very long strings that might
    exceed system limits,
2 -- be truncated, or be unwieldy for display or parsing. It also
    denormalizes data into a string,
3 -- making SQL operations on individual elements (e.g., finding
    departments with an employee
4 -- named 'Alice' within the aggregated string) complex and inefficient.
5
6 -- Alternative for relational processing:
7 SELECT
8     COALESCE(d.departmentName, 'No Department Assigned') AS
    departmentName,
9     e.firstName
10 FROM Employees e
11 LEFT JOIN Departments d ON e.departmentId = d.departmentId
12 ORDER BY COALESCE(d.departmentName, 'No Department Assigned') NULLS
    FIRST, e.firstName;
```

Listing 3: Solution 1.1.2: Disadvantage of STRING_AGG and alternative

1.1.3 Practice Inefficient Alternatives Avoidance

```
1 -- Inefficient: Fetching all skills for Engineering employees and then
    using application code
2 -- (e.g., Python loop) to build the unique, sorted list.
3
4 -- Efficient SQL approach using STRING_AGG with UNNEST for array '
    skills' column:
5 SELECT
6     d.departmentName,
7     STRING_AGG(DISTINCT skill, '; ' ORDER BY skill) AS
    departmentUniqueSkills
8 FROM Employees e
9 JOIN Departments d ON e.departmentId = d.departmentId,
10 UNNEST(e.skills) AS skill -- Creates a row for each skill in the array
11 WHERE d.departmentName = 'Engineering'
```

```
12 GROUP BY d.departmentId, d.departmentName;
```

Listing 4: Solution 1.1.3: Efficient unique skills list

1.2 ARRAY_AGG(expression [ORDER BY ...])

1.2.1 Practice Meaning, Values, Relations, Advantages

```
1 SELECT
2     p.projectName,
3     ARRAY_AGG(ep.employeeId ORDER BY ep.employeeId) AS teamMemberIds
4 FROM Projects p
5 JOIN EmployeeProjects ep ON p.projectId = ep.projectId
6 GROUP BY p.projectId, p.projectName
7 ORDER BY p.projectName;
```

Listing 5: Solution 1.2.1: Project team member IDs array

1.2.2 Practice Disadvantages

```
1 -- Disadvantage: Querying specific positions or complex conditions
   -- within arrays
2 -- (e.g., "employeeId X is the first element and employeeId Y is the
   -- third")
3 -- can be cumbersome and less performant in SQL compared to querying
   -- normalized data.
4 -- Standard SQL indexes are generally not effective for searching
   -- inside array elements
5 -- by position or value across many rows. While PostgreSQL has good
   -- array functions,
6 -- complex array logic can be slower.
7
8 -- Normalized structure (EmployeeProjects table) is better for such
   -- specific relational queries:
9 SELECT ep.projectId, ep.employeeId -- (plus ordering or window
   -- functions if "first assigned" means based on a date or sequence)
10 FROM EmployeeProjects ep
11 ORDER BY ep.projectId, ep.assignmentId; -- Assuming assignmentId
   -- implies order
12
13 -- To find if employeeId 101 is the first in an array created by
   -- ARRAY_AGG(ep.employeeId ORDER BY ep.assignmentId):
14 -- SELECT p.projectName
15 -- FROM Projects p
16 -- JOIN (
17 --     SELECT
18 --         projectId,
19 --         ARRAY_AGG(employeeId ORDER BY assignmentId) as ids
20 --     FROM EmployeeProjects
21 --     GROUP BY projectId
22 -- ) AS pa ON p.projectId = pa.projectId
23 -- WHERE pa.ids[1] = 101;
24 -- This is possible but often less direct than querying normalized
   -- tables with appropriate indexing.
```

Listing 6: Solution 1.2.2: Disadvantage of ARRAY_AGG for positional queries

1.2.3 Practice Inefficient Alternatives Avoidance

```
1 -- Inefficient: N+1 queries (1 query for categories, then N queries for
  products in each category).
2
3 -- Efficient SQL approach:
4 SELECT
5     COALESCE(p.category, 'Uncategorized') AS productCategory,
6     ARRAY_AGG(p.productName ORDER BY p.productName) AS
  productNamesInCategory
7 FROM Products p
8 GROUP BY p.category
9 ORDER BY productCategory;
```

Listing 7: Solution 1.2.3: Efficient product names per category

1.3 JSON_AGG(expression [ORDER BY ...])

1.3.1 Practice Meaning, Values, Relations, Advantages

```
1 SELECT
2     d.departmentName ,
3     JSON_AGG(
4         JSON_BUILD_OBJECT(
5             'firstName', e.firstName ,
6             'lastName', e.lastName ,
7             'salary', e.salary
8         )
9         ORDER BY e.salary DESC
10    ) AS employeesJson
11 FROM Departments d
12 JOIN Employees e ON d.departmentId = e.departmentId
13 WHERE d.locationCity = 'San Francisco'
14 GROUP BY d.departmentId, d.departmentName
15 ORDER BY d.departmentName;
```

Listing 8: Solution 1.3.1: JSON array of employees in San Francisco departments

1.3.2 Practice Disadvantages

```
1 -- Disadvantage: Generating very large JSON documents (e.g., thousands
  of employees per department,
2 -- each with many fields) can consume significant memory and CPU on the
  database server.
3 -- The resulting JSON string can also be large to transmit and parse by
  the client.
4 -- JSON is inherently schema-less; while JSON_BUILD_OBJECT helps
  structure it,
5 -- the database doesn't enforce the internal schema of the JSON object
  in the same way
6 -- as table columns. Consumers must implement their own validation or
  be robust to variations
7 -- if the JSON structure changes.
```



```

8 SELECT 'Large JSON objects from JSON_AGG can strain server resources
   and network. Type safety relies on producer/consumer discipline.' AS
   DisadvantageInfo;

```

Listing 9: Solution 1.3.2: Disadvantage of JSON_AGG for large/complex objects

1.3.3 Practice Inefficient Alternatives Avoidance

```

1  -- Inefficient: Multiple database queries (N+1 problem) and manual JSON
   construction in application code.
2
3  -- Efficient SQL approach:
4  SELECT
5      p.productName,
6      JSON_AGG(
7          JSON_BUILD_OBJECT(
8              'saleId', s.saleId,
9              'saleDate', s.saleDate,
10             'quantity', s.quantity,
11             'regionName', rg.regionName,
12             'saleNotes', s.notes
13         )
14         ORDER BY s.saleDate
15     ) FILTER (WHERE s.saleId IS NOT NULL) AS salesJson -- FILTER
   ensures empty JSON array [] if no sales, not [null]
16 FROM Products p
17 LEFT JOIN Sales s ON p.productId = s.productId
18 LEFT JOIN Regions rg ON s.regionId = rg.regionId
19 GROUP BY p.productId, p.productName
20 ORDER BY p.productName;

```

Listing 10: Solution 1.3.3: Efficient JSON feed of products and sales

1.4 PERCENTILE_CONT(fraction) WITHIN GROUP (ORDER BY sort_expression)

1.4.1 Practice Meaning, Values, Relations, Advantages

```

1 SELECT
2     p.category AS productCategory,
3     PERCENTILE_CONT(0.25) WITHIN GROUP (ORDER BY p.listPrice) AS
   p25ListPrice,
4     PERCENTILE_CONT(0.50) WITHIN GROUP (ORDER BY p.listPrice) AS
   medianListPrice,
5     PERCENTILE_CONT(0.75) WITHIN GROUP (ORDER BY p.listPrice) AS
   p75ListPrice
6 FROM Products p
7 WHERE p.category IS NOT NULL
8 GROUP BY p.category
9 ORDER BY p.category;

```

Listing 11: Solution 1.4.1: Percentiles of listPrice per category

1.4.2 Practice Disadvantages

```

1  -- Disadvantage: PERCENTILE_CONT interpolates between values to find
    the percentile,
2  -- which can result in a value not actually present in the dataset.
3  -- For example, if values are 1 and 2, the 0.5 percentile (median) is
    1.5.
4  -- If the data is inherently discrete (e.g., star ratings 1-5), and you
    need a percentile value
5  -- that IS one of the actual data points, PERCENTILE_DISC might be more
    appropriate
6  -- as it picks an existing value.
7  -- The choice depends on whether a continuous or discrete
    interpretation of percentiles is desired.
8  SELECT 'PERCENTILE_CONT interpolates, potentially yielding values not
    in the dataset. For discrete data, PERCENTILE_DISC might be more
    intuitive if an existing value is required.' AS DisadvantageInfo;

```

Listing 12: Solution 1.4.2: Disadvantage of PERCENTILE_CONT interpolation

1.4.3 Practice Inefficient Alternatives Avoidance

```

1  -- Inefficient: Data export and manual/spreadsheet calculation for a
    standard statistical measure.
2
3  -- Efficient SQL approach:
4  SELECT
5      COALESCE(d.departmentName, 'No Department Assigned') AS
        departmentName,
6      PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY e.salary) AS
        medianSalary
7  FROM Employees e
8  LEFT JOIN Departments d ON e.departmentId = d.departmentId
9  GROUP BY d.departmentId, d.departmentName
10 ORDER BY medianSalary DESC NULLS LAST;

```

Listing 13: Solution 1.4.3: Efficient median salary calculation

1.5 CORR(Y, X)

1.5.1 Practice Meaning, Values, Relations, Advantages

```

1  SELECT
2      CORR(s.quantity, p.listPrice) AS quantityPriceCorrelation
3  FROM Sales s
4  JOIN Products p ON s.productId = p.productId;

```

Listing 14: Solution 1.5.1: Correlation between sales quantity and list price

1.5.2 Practice Disadvantages

```

1  -- Disadvantage: CORR only measures linear association. A correlation
    coefficient near 0
2  -- means there is little to no linear relationship.
3  -- However, it does NOT mean there is no relationship at all. There
    could be a strong
4  -- non-linear relationship (e.g., quadratic, U-shaped) for which CORR
    would be close to 0.

```

```

5 -- Example: If sales quantity is high for very cheap and very expensive
   products but low
6 -- for mid-priced ones (a U-shape), CORR(quantity, price) might be near
   0.
7 SELECT 'CORR only captures linear relationships. A value near 0 does
   not rule out strong non-linear relationships.' AS DisadvantageInfo;

```

Listing 15: Solution 1.5.2: Disadvantage of CORR for non-linear relationships

1.5.3 Practice Inefficient Alternatives Avoidance

```

1 -- Inefficient: Exporting data to external software for a basic
   correlation calculation.
2
3 -- Efficient SQL approach:
4 SELECT
5     CORR(e.performanceScore, e.salary) AS salaryPerformanceCorrelation
6 FROM Employees e
7 WHERE e.performanceScore IS NOT NULL AND e.salary IS NOT NULL; -- CORR
   ignores NULL pairs

```

Listing 16: Solution 1.5.3: Efficient salary-performance correlation

1.6 REGR_SLOPE(Y, X)

1.6.1 Practice Meaning, Values, Relations, Advantages

```

1 SELECT
2     p.category,
3     REGR_SLOPE(s.quantity, p.listPrice) AS salesQtyVsPriceSlope
4 FROM Sales s
5 JOIN Products p ON s.productId = p.productId
6 WHERE p.category = 'Electronics'
7 GROUP BY p.category;
8 -- This gives the slope for individual sales transactions. For slope of
   average quantity
9 -- vs price, a subquery to average quantity per product first would be
   needed.

```

Listing 17: Solution 1.6.1: Regression slope for Electronics sales quantity vs. price

1.6.2 Practice Disadvantages

```

1 -- Disadvantage: REGR_SLOPE only provides the slope of the linear
   regression line. It does not indicate:
2 -- 1. The strength of the linear relationship (e.g., R-squared value,
   provided by REGR_R2).
3 -- 2. A slope might be calculated even if data points are widely
   scattered and the linear model is a poor fit.
4 -- 3. Whether a linear model is appropriate at all (data could be non-
   linear).
5 -- 4. The statistical significance of the slope (p-value).
6 SELECT 'REGR_SLOPE does not indicate goodness-of-fit (like R-squared)
   or statistical significance of the relationship.' AS
   DisadvantageInfo;

```

Listing 18: Solution 1.6.2: Disadvantage of REGR_SLOPE regarding fit/significance

1.6.3 Practice Inefficient Alternatives Avoidance

```
1  -- Inefficient: Exporting data to Excel for plotting and trendline
   slope calculation for a simple linear regression.
2
3  -- Efficient SQL approach:
4  SELECT
5      d.departmentName ,
6      REGR_SLOPE(e.performanceScore , e.salary) AS perfScoreVsSalarySlope
7  FROM Employees e
8  JOIN Departments d ON e.departmentId = d.departmentId
9  WHERE d.departmentName = 'Sales' AND e.performanceScore IS NOT NULL
10 GROUP BY d.departmentId, d.departmentName;
```

Listing 19: Solution 1.6.3: Efficient performance vs. salary slope for Sales dept

2 Category: Advanced Grouping Operations

2.1 GROUPING SETS ((set1), (set2), ...)

2.1.1 Practice Meaning, Values, Relations, Advantages

```
1 SELECT
2     COALESCE(p.category, 'All Categories') AS productCategory,
3     COALESCE(r.regionName, 'All Regions') AS region,
4     SUM(s.quantity) AS totalQuantitySold,
5     SUM(p.listPrice * s.quantity) AS totalListPriceValue,
6     GROUPING(p.category, r.regionName) AS groupingIndicator
7 FROM Sales s
8 JOIN Products p ON s.productId = p.productId
9 LEFT JOIN Regions r ON s.regionId = r.regionId
10 GROUP BY GROUPING SETS (
11     (p.category, r.regionName),
12     (p.category),
13     (r.regionName),
14     ()
15 )
16 ORDER BY GROUPING(p.category, r.regionName), productCategory NULLS LAST
17         , region NULLS LAST;
```

Listing 20: Solution 2.1.1: Sales summary with multiple GROUPING SETS

2.1.2 Practice Disadvantages

```
1 -- Disadvantage: With many columns and/or many specific grouping sets,
2   the GROUPING SETS
3   clause can become very long and complex to write correctly.
4 -- It increases the risk of errors in defining the desired combinations
5   , or omitting a
6   necessary one / including an unnecessary one.
7 -- The resulting output can also be large and harder to interpret if
8   not carefully
9   structured and labelled using GROUPING() and COALESCE.
10 -- Debugging such complex queries can be challenging.
11 SELECT 'Highly complex GROUPING SETS definitions can be error-prone to
12 write and maintain, and their output hard to manage.' AS
13 DisadvantageInfo;
```

Listing 21: Solution 2.1.2: Disadvantage of complex GROUPING SETS

2.1.3 Practice Inefficient Alternatives Avoidance

```
1 -- Inefficient/Basic (conceptual):
2 /*
3 SELECT EXTRACT(YEAR FROM s.saleDate) AS saleYear, p.category, SUM(s.
4 quantity) AS totalQuantity
5 FROM Sales s JOIN Products p ON s.productId = p.productId
6 GROUP BY EXTRACT(YEAR FROM s.saleDate), p.category
7 UNION ALL
8 SELECT EXTRACT(YEAR FROM s.saleDate) AS saleYear, NULL AS category, SUM
9 (s.quantity) AS totalQuantity
```

```

8 FROM Sales s JOIN Products p ON s.productId = p.productId -- Join might
   not be needed if category not selected
9 GROUP BY EXTRACT(YEAR FROM s.saleDate);
10 */
11
12 -- Efficient GROUPING SETS approach:
13 SELECT
14     EXTRACT(YEAR FROM s.saleDate) AS saleYear,
15     COALESCE(p.category, 'All Categories for Year') AS productCategory,
16     SUM(s.quantity) AS totalQuantity
17 FROM Sales s
18 JOIN Products p ON s.productId = p.productId
19 GROUP BY GROUPING SETS (
20     (EXTRACT(YEAR FROM s.saleDate), p.category),
21     (EXTRACT(YEAR FROM s.saleDate))
22 )
23 ORDER BY saleYear, productCategory NULLS LAST;

```

Listing 22: Solution 2.1.3: Efficient sales quantity by year/category with GROUPING SETS

2.2 ROLLUP (col1, col2, ...)

2.2.1 Practice Meaning, Values, Relations, Advantages

```

1 SELECT
2     COALESCE(d.departmentName, 'Overall Total') AS department,
3     CASE WHEN GROUPING(d.departmentName) = 0 THEN COALESCE(p.
   projectName, 'Department Total') ELSE NULL END AS project,
4     SUM(ep.hoursWorked) AS totalHours,
5     GROUPING(d.departmentName, p.projectName) AS groupingLevel
6 FROM EmployeeProjects ep
7 JOIN Employees e ON ep.employeeId = e.employeeId
8 LEFT JOIN Departments d ON e.departmentId = d.departmentId
9 JOIN Projects p ON ep.projectId = p.projectId
10 GROUP BY ROLLUP (d.departmentName, p.projectName)
11 ORDER BY d.departmentName NULLS LAST, p.projectName NULLS LAST;

```

Listing 23: Solution 2.2.1: Hierarchical summary of hours worked with ROLLUP

2.2.2 Practice Disadvantages

```

1 -- Disadvantage: ROLLUP is designed for strict hierarchical
   summarization. The order of
2 -- columns is critical and it only "rolls up" from right to left.
3 -- ROLLUP(country, state, city) will NOT produce a subtotal for (
   country, city) by itself
4 -- (skipping state) or a subtotal for (city) by itself.
5 -- To get such non-standard hierarchical or cross-level subtotals, you
   would need to use
6 -- GROUPING SETS to specify those exact combinations.
7 -- Relying solely on ROLLUP limits you to its predefined hierarchical
   aggregation paths.

```

```

8 SELECT 'ROLLUP is inflexible for non-hierarchical subtotals (e.g., (
   country, city) from ROLLUP(country, state, city)). Use GROUPING SETS
   for such cases.' AS DisadvantageInfo;

```

Listing 24: Solution 2.2.2: Disadvantage of ROLLUP for non-hierarchical subtotals

2.2.3 Practice Inefficient Alternatives Avoidance

```

1 -- Inefficient: Multiple UNION ALL queries to build up hierarchical
   subtotals.
2
3 -- Efficient ROLLUP approach:
4 SELECT
5     COALESCE(rg.regionName, 'Grand Total') AS region,
6     CASE WHEN GROUPING(rg.regionName) = 0 THEN COALESCE(p.category, '
   Region Total') ELSE NULL END AS category,
7     CASE WHEN GROUPING(rg.regionName, p.category) = 0 THEN COALESCE(p.
   productName, 'Category Total') ELSE NULL END AS productName,
8     SUM(s.quantity) AS totalQuantity
9 FROM Sales s
10 JOIN Products p ON s.productId = p.productId
11 LEFT JOIN Regions rg ON s.regionId = rg.regionId
12 GROUP BY ROLLUP (rg.regionName, p.category, p.productName)
13 ORDER BY rg.regionName NULLS LAST, p.category NULLS LAST, p.productName
   NULLS LAST;

```

Listing 25: Solution 2.2.3: Efficient hierarchical sales report with ROLLUP

2.3 CUBE (col1, col2, ...)

2.3.1 Practice Meaning, Values, Relations, Advantages

```

1 SELECT
2     EXTRACT(YEAR FROM s.saleDate) AS saleYear,
3     COALESCE(p.category, 'All Categories') AS productCategory,
4     SUM(s.quantity) AS totalQuantity,
5     GROUPING(EXTRACT(YEAR FROM s.saleDate), p.category) AS
   groupingIndicator
6 FROM Sales s
7 JOIN Products p ON s.productId = p.productId
8 GROUP BY CUBE (EXTRACT(YEAR FROM s.saleDate), p.category)
9 ORDER BY saleYear NULLS LAST, productCategory NULLS LAST;

```

Listing 26: Solution 2.3.1: Cross-tabular sales summary with CUBE

2.3.2 Practice Disadvantages

```

1 -- Disadvantage: CUBE generates all possible combinations of subtotals
   (2^N for N columns).
2 -- If N is large, this results in a very large number of rows in the
   output, many of
3 -- which may not be relevant to the user.
4 -- This can lead to:
5 -- 1. Increased query processing time and resource consumption on the
   database.

```

```

6 -- 2. Large result sets that are difficult for users to navigate or for
   applications to process.
7 -- It's often better to use GROUPING SETS to specify only the required
   combinations if not all 2^N are needed.
8 SELECT 'CUBE can produce excessive, unneeded subtotals for many columns
   , leading to performance issues and unwieldy results.' AS
   DisadvantageInfo;

```

Listing 27: Solution 2.3.2: Disadvantage of CUBE generating excessive subtotals

2.3.3 Practice Inefficient Alternatives Avoidance

```

1 -- Inefficient: Running multiple separate queries or a complex UNION
   ALL of queries for each desired aggregation level.
2
3 -- Efficient CUBE approach:
4 SELECT
5     COALESCE(r.regionName, 'All Regions') AS region,
6     COALESCE(p.category, 'All Categories') AS productCategory,
7     SUM(s.quantity) AS totalQuantity
8 FROM Sales s
9 LEFT JOIN Regions r ON s.regionId = r.regionId
10 JOIN Products p ON s.productId = p.productId
11 GROUP BY CUBE (r.regionName, p.category)
12 ORDER BY r.regionName NULLS LAST, p.category NULLS LAST;

```

Listing 28: Solution 2.3.3: Efficient multi-level sales exploration with CUBE

3 Hardcore Combined Problem

```
1 WITH EmployeeBasicInfo AS (  
2     SELECT  
3         e.employeeId,  
4         e.firstName || ' ' || e.lastName AS employeeFullName,  
5         EXTRACT(YEAR FROM e.hireDate) AS employeeHireYear,  
6         e.salary,  
7         e.departmentId,  
8         e.performanceScore,  
9         COALESCE(STRING_AGG(skill, ', ' ORDER BY skill), 'No Skills  
10        Listed') AS skillsList  
11     FROM Employees e  
12     LEFT JOIN UNNEST(e.skills) AS skill ON TRUE -- UNNEST skills array  
13     GROUP BY e.employeeId, e.firstName, e.lastName, e.hireDate, e.  
14         salary, e.departmentId, e.performanceScore  
15 ),  
16 EmployeeProjectSummary AS (  
17     SELECT  
18         ep.employeeId,  
19         COALESCE(SUM(ep.hoursWorked), 0) AS totalHoursOnProjects,  
20         COALESCE(STRING_AGG(DISTINCT p.projectName, ', ' ORDER BY p.  
21         projectName), 'None') AS projectsParticipated  
22     FROM EmployeeProjects ep  
23     JOIN Projects p ON ep.projectId = p.projectId  
24     GROUP BY ep.employeeId  
25 ),  
26 SalesWithSatisfaction2023 AS (  
27     SELECT  
28         s.employeeId,  
29         s.quantity,  
30         CAST(NULLIF(s.notes ->> 'customerSatisfaction', '') AS INTEGER)  
31         AS satisfactionScore  
32     FROM Sales s  
33     WHERE EXTRACT(YEAR FROM s.saleDate) = 2023  
34         AND s.employeeId IS NOT NULL  
35         AND s.notes ->> 'customerSatisfaction' IS NOT NULL  
36         AND s.notes ->> 'customerSatisfaction' ~ '^[0-9]+$' -- Ensures  
37         the string is purely numeric  
38 ),  
39 EmployeeCorrelation2023 AS (  
40     SELECT  
41         employeeId,  
42         CORR(satisfactionScore, quantity) AS  
43         salesQtyVsSatisfactionCorr2023  
44     FROM SalesWithSatisfaction2023  
45     GROUP BY employeeId  
46     HAVING COUNT(satisfactionScore) >= 2 -- CORR needs at least 2 pairs  
47         of non-null values  
48 ),  
49 EmployeeSalesRevenue2023 AS (  
50     SELECT  
51         s.employeeId,  
52         COALESCE(SUM(s.quantity * pr.listPrice), 0.00) AS  
53         totalRevenueGenerated2023  
54     FROM Sales s  
55     JOIN Products pr ON s.productId = pr.productId  
56     WHERE EXTRACT(YEAR FROM s.saleDate) = 2023 AND s.employeeId IS NOT
```

```

49     NULL
50     GROUP BY s.employeeId
51 ),
52 DepartmentalAggregates AS (
53     SELECT
54         e.departmentId, -- This will be NULL for employees without a
55         department
56         PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY e.salary) AS
57         medianSalaryInDepartment,
58         JSON_AGG(
59             JSON_BUILD_OBJECT(
60                 'employeeId', e.employeeId,
61                 'performanceScore', e.performanceScore
62             )
63             ORDER BY e.performanceScore DESC
64         ) FILTER (WHERE e.performanceScore IS NOT NULL) AS
65         departmentPerformanceOverviewJson
66     FROM Employees e
67     GROUP BY e.departmentId -- Handles NULL departmentId as a separate
68     group
69 )
70 SELECT
71     CASE
72         WHEN GROUPING(COALESCE(dpt.departmentName, 'No Department
73         Assigned')), ebi.employeeId) = 0 THEN 'Employee Detail'
74         WHEN GROUPING(ebi.employeeId) = 1 AND GROUPING(COALESCE(dpt.
75         departmentName, 'No Department Assigned')) = 0 THEN 'Department
76         Summary'
77         ELSE 'Grand Total'
78     END AS reportingLevel,
79     CASE
80         WHEN GROUPING(COALESCE(dpt.departmentName, 'No Department
81         Assigned')) = 1 THEN 'Overall Summary' -- For Grand Total row
82         ELSE COALESCE(dpt.departmentName, 'No Department Assigned')
83     END AS departmentName,
84     CASE WHEN GROUPING(ebi.employeeId) = 0 THEN ebi.employeeFullName
85     ELSE NULL END AS employeeFullName,
86     CASE WHEN GROUPING(ebi.employeeId) = 0 THEN ebi.employeeHireYear
87     ELSE NULL END AS employeeHireYear,
88     CASE WHEN GROUPING(ebi.employeeId) = 0 THEN ebi.skillsList ELSE
89     NULL END AS skillsList,
90     CASE WHEN GROUPING(ebi.employeeId) = 0 THEN COALESCE(eps.
91     projectsParticipated, 'None') ELSE NULL END AS projectsParticipated,
92     SUM(COALESCE(eps.totalHoursOnProjects, 0)) AS totalHoursOnProjects,
93     SUM(COALESCE(esr.totalRevenueGenerated2023, 0.00)) AS
94     totalRevenueGenerated2023,
95     CASE WHEN GROUPING(ebi.employeeId) = 0 THEN ec.
96     salesQtyVsSatisfactionCorr2023 ELSE NULL END AS
97     salesQtyVsSatisfactionCorr2023,
98     CASE
99         WHEN GROUPING(ebi.employeeId) = 1 AND GROUPING(COALESCE(dpt.
100         departmentName, 'No Department Assigned')) = 0 THEN da.
101         medianSalaryInDepartment

```

```

89         ELSE NULL
90     END AS medianSalaryInDepartment,
91     CASE
92         WHEN GROUPING(ebi.employeeId) = 1 AND GROUPING(COALESCE(dpt.
departmentName, 'No Department Assigned')) = 0 THEN COALESCE(da.
departmentPerformanceOverviewJson, '[]'::JSONB)
93         ELSE NULL
94     END AS departmentPerformanceOverviewJson
95
96 FROM EmployeeBasicInfo ebi
97 LEFT JOIN Departments dpt ON ebi.departmentId = dpt.departmentId
98 LEFT JOIN EmployeeProjectSummary eps ON ebi.employeeId = eps.employeeId
99 LEFT JOIN EmployeeSalesRevenue2023 esr ON ebi.employeeId = esr.
employeeId
100 LEFT JOIN EmployeeCorrelation2023 ec ON ebi.employeeId = ec.employeeId
101 LEFT JOIN DepartmentalAggregates da ON ebi.departmentId IS NOT DISTINCT
FROM da.departmentId -- Handles NULL departmentId for 'No
Department Assigned' group
102
103 GROUP BY GROUPING SETS (
104     -- Employee Detail Level (includes department-level aggregates that
will be picked by CASE based on grouping)
105     (COALESCE(dpt.departmentName, 'No Department Assigned'), ebi.
employeeId, ebi.employeeFullName, ebi.employeeHireYear, ebi.
skillsList, eps.projectsParticipated, ec.
salesQtyVsSatisfactionCorr2023, da.medianSalaryInDepartment, da.
departmentPerformanceOverviewJson),
106     -- Department Summary Level
107     (COALESCE(dpt.departmentName, 'No Department Assigned'), da.
medianSalaryInDepartment, da.departmentPerformanceOverviewJson),
108     -- Grand Total Level
109     ()
110 )
111 ORDER BY
112     CASE
113         WHEN GROUPING(COALESCE(dpt.departmentName, 'No Department
Assigned')) = 1 THEN 3 -- Grand Total last
114         WHEN dpt.departmentName IS NULL AND COALESCE(dpt.departmentName
, 'No Department Assigned') = 'No Department Assigned' THEN 1 -- 'No
Department Assigned' first
115         ELSE 2 -- Actual departments
116     END,
117     departmentName,
118     reportingLevel,
119     employeeFullName NULLS LAST;

```

Listing 29: Solution to Hardcore Combined Problem