

Aggregate Functions in SQL: Exercises

Generated Exercises

May 8, 2025

Contents

1	Practice Meanings, Values, Relations, and Advantages	5
2	Practice Disadvantages and Potential Pitfalls	5
3	Inefficient vs. Efficient Aggregate Usage	6
4	Hardcore Problem Combining Concepts	7

Dataset for PostgreSQL

The following SQL code creates and populates the necessary tables for the exercises.

```
1  -- Dataset for Aggregate Functions Exercises (PostgreSQL)
2
3  -- Drop tables if they exist to ensure a clean setup
4  DROP TABLE IF EXISTS employee_tasks CASCADE;
5  DROP TABLE IF EXISTS projects CASCADE;
6  DROP TABLE IF EXISTS employees CASCADE;
7  DROP TABLE IF EXISTS departments CASCADE;
8
9  -- departments table
10 CREATE TABLE departments (
11     department_id SERIAL PRIMARY KEY,
12     department_name VARCHAR(100) NOT NULL UNIQUE,
13     location VARCHAR(100)
14 );
15
16 -- employees table
17 CREATE TABLE employees (
18     employee_id SERIAL PRIMARY KEY,
19     first_name VARCHAR(50) NOT NULL,
20     last_name VARCHAR(50) NOT NULL,
21     email VARCHAR(100) UNIQUE,
22     salary DECIMAL(10, 2) CHECK (salary > 0),
23     hire_date DATE,
24     department_id INT REFERENCES departments(department_id),
25     manager_id INT REFERENCES employees(employee_id),
26     bonus_percentage NUMERIC(4, 2) CHECK (bonus_percentage >= 0 AND bonus_percentage <= 100.00),
27     performance_rating INT CHECK (performance_rating >= 1 AND performance_rating <= 5)
28 );
29
30 -- projects table
31 CREATE TABLE projects (
32     project_id SERIAL PRIMARY KEY,
33     project_name VARCHAR(150) NOT NULL UNIQUE,
34     start_date DATE,
35     end_date DATE,
36     budget DECIMAL(12, 2) CHECK (budget > 0),
37     lead_employee_id INT REFERENCES employees(employee_id)
38 );
39
40 -- employee_tasks table
41 CREATE TABLE employee_tasks (
42     task_id SERIAL PRIMARY KEY,
43     employee_id INT NOT NULL REFERENCES employees(employee_id),
44     project_id INT NOT NULL REFERENCES projects(project_id),
45     task_description TEXT,
46     hours_spent DECIMAL(5, 2) CHECK (hours_spent >= 0),
47     task_date DATE,
48     status VARCHAR(20) DEFAULT 'Pending' CHECK (status IN ('Pending', 'In Progress', 'Completed', 'Cancelled'))
49 );
50
51 -- Populate departments
52 INSERT INTO departments (department_name, location) VALUES
53 ('Human Resources', 'New York'),
54 ('Engineering', 'San Francisco'),
55 ('Marketing', 'Chicago'),
56 ('Sales', 'Boston'),
57 ('Research', 'Austin');
58
59 -- Populate employees
60 INSERT INTO employees (first_name, last_name, email, salary, hire_date, department_id,
61     manager_id, bonus_percentage, performance_rating) VALUES
62 ('Alice', 'Smith', 'alice.smith@example.com', 70000.00, '2020-01-15', 2, NULL, 10.00, 4),
63     -- Manager for Eng
64 ('Bob', 'Johnson', 'bob.johnson@example.com', 120000.00, '2019-03-01', 2, 1, 15.00, 5),
65 ('Charlie', 'Brown', 'charlie.brown@example.com', 95000.00, '2021-07-20', 2, 1, 12.50, 4),
66 ('Diana', 'Lee', 'diana.lee@example.com', 80000.00, '2022-06-10', 2, 1, NULL, 3),
```

```

65 ('Eve', 'Davis', 'eve.davis@example.com', 60000.00, '2020-05-01', 1, NULL, 5.00, 4), --
    Manager for HR
66 ('Frank', 'Miller', 'frank.miller@example.com', 55000.00, '2021-08-25', 1, 5, 7.50, 3),
67 ('Grace', 'Wilson', 'grace.wilson@example.com', 58000.00, '2022-01-10', 1, 5, NULL, 4),
68 ('Henry', 'Moore', 'henry.moore@example.com', 90000.00, '2019-11-05', 3, NULL, 10.00, 5)
    , -- Manager for Marketing
69 ('Ivy', 'Taylor', 'ivy.taylor@example.com', 75000.00, '2020-02-17', 3, 8, 8.00, 4),
70 ('Jack', 'Anderson', 'jack.anderson@example.com', 72000.00, '2021-09-30', 3, 8, 9.50, 3)
    ,
71 ('Karen', 'Thomas', 'karen.thomas@example.com', 65000.00, '2023-01-20', 3, 8, NULL, 4),
72 ('Leo', 'Jackson', 'leo.jackson@example.com', 110000.00, '2018-07-14', 4, NULL, 20.00,
    5), -- Manager for Sales
73 ('Mia', 'White', 'mia.white@example.com', 85000.00, '2019-04-01', 4, 12, 18.00, 4),
74 ('Noah', 'Harris', 'noah.harris@example.com', 82000.00, '2020-10-15', 4, 12, 17.50, 4),
75 ('Olivia', 'Martin', 'olivia.martin@example.com', 78000.00, '2021-12-01', 4, 12, NULL,
    3),
76 ('Paul', 'Garcia', 'paul.garcia@example.com', 130000.00, '2017-05-22', 5, NULL, 22.00,
    5), -- Manager for Research
77 ('Quinn', 'Martinez', 'quinn.martinez@example.com', 100000.00, '2018-09-10', 5, 16,
    15.00, 5),
78 ('Ruby', 'Robinson', 'ruby.robinson@example.com', 92000.00, '2019-06-05', 5, 16, 14.00,
    4),
79 ('Sam', 'Clark', 'sam.clark@example.com', 88000.00, '2020-11-11', 5, 16, NULL, 3),
80 ('Tina', 'Rodriguez', 'tina.rodriguez@example.com', 105000.00, '2022-03-15', 2, 1,
    10.00, 5),
81 ('Uma', 'Lewis', 'uma.lewis@example.com', 62000.00, '2021-05-10', 1, 5, 6.00, 4),
82 ('Victor', 'Walker', 'victor.walker@example.com', 77000.00, '2022-08-01', 3, 8, 7.00, 3)
    ,
83 ('Wendy', 'Hall', 'wendy.hall@example.com', 90000.00, '2020-01-20', 4, 12, 19.00, 5),
84 ('Xavier', 'Allen', 'xavier.allen@example.com', 115000.00, '2021-02-18', 5, 16, 16.00,
    4),
85 ('Yara', 'Young', 'yara.young@example.com', 71000.00, '2023-04-01', 2, 1, 5.00, 3);
86
87 -- Update manager_id for Alice, Eve, Henry, Leo, Paul (they were NULL, now they are
    their own managers for simplicity or a designated top manager if exists)
88 -- For this dataset, let's assume they are top-level or report to someone outside this
    scope.
89 -- Or, let's make Alice the overall CEO reporting to no one.
90 UPDATE employees SET manager_id = 1 WHERE employee_id IN (5, 8, 12, 16); -- Other
    managers report to Alice
91 UPDATE employees SET manager_id = NULL WHERE employee_id = 1;
92
93
94 -- Populate projects
95 INSERT INTO projects (project_name, start_date, end_date, budget, lead_employee_id)
    VALUES
96 ('Project Alpha', '2023-01-15', '2023-06-30', 150000.00, 2), -- Bob (Eng)
97 ('Project Beta', '2023-03-01', '2023-09-30', 250000.00, 17), -- Quinn (Research)
98 ('Project Gamma', '2023-05-10', '2023-12-31', 100000.00, 9), -- Ivy (Marketing)
99 ('Project Delta', '2023-07-01', '2024-01-31', 300000.00, 13), -- Mia (Sales)
100 ('Project Epsilon', '2023-09-01', '2024-03-31', 75000.00, 3), -- Charlie (Eng)
101 ('Project Zeta', '2024-01-01', '2024-06-30', 220000.00, 18); -- Ruby (Research)
102
103 -- Populate employee_tasks
104 INSERT INTO employee_tasks (employee_id, project_id, task_description, hours_spent,
    task_date, status) VALUES
105 (2, 1, 'Initial design phase', 40.5, '2023-01-20', 'Completed'),
106 (3, 1, 'Frontend development', 120.0, '2023-03-15', 'Completed'),
107 (4, 1, 'Backend development', 150.75, '2023-04-10', 'In Progress'),
108 (20, 1, 'API integration', 80.0, '2023-05-01', 'In Progress'),
109 (17, 2, 'Literature review', 60.0, '2023-03-10', 'Completed'),
110 (18, 2, 'Experiment setup', 90.5, '2023-05-01', 'In Progress'),
111 (19, 2, 'Data collection', 70.0, '2023-06-15', 'Pending'),
112 (24, 2, 'Preliminary analysis', 30.25, '2023-07-01', 'Pending'),
113 (9, 3, 'Market research', 50.0, '2023-05-15', 'Completed'),
114 (10, 3, 'Campaign strategy', 75.25, '2023-06-20', 'In Progress'),
115 (11, 3, 'Content creation', 100.0, '2023-08-01', 'In Progress'),
116 (22, 3, 'Social media outreach', 40.0, '2023-07-10', 'Pending'),
117 (13, 4, 'Lead generation plan', 45.0, '2023-07-05', 'Completed'),
118 (14, 4, 'Client outreach', 110.0, '2023-08-15', 'In Progress'),
119 (15, 4, 'Sales calls', 95.5, '2023-09-01', 'In Progress'),
120 (23, 4, 'Contract negotiation', 60.75, '2023-09-20', 'Pending'),
121 (3, 5, 'Requirements gathering', 30.0, '2023-09-05', 'Completed'),

```

```

122 (4, 5, 'Prototyping', 80.0, '2023-10-10', 'In Progress'),
123 (20, 5, 'User testing setup', 40.5, '2023-11-01', 'Pending'),
124 (25, 5, 'Documentation', 25.0, '2023-11-15', 'Pending'),
125 (18, 6, 'Advanced algorithm design', 120.0, '2024-01-10', 'In Progress'),
126 (19, 6, 'Simulation runs', 100.0, '2024-02-15', 'Pending'),
127 (24, 6, 'Results validation', 70.25, '2024-03-01', 'Pending'),
128 (6, 1, 'HR support for Alpha team', 10.0, '2023-02-01', 'Completed'),
129 (7, 3, 'HR support for Gamma team', 12.0, '2023-06-01', 'Completed'),
130 (2, 1, 'Additional design review', 15.0, '2023-02-20', 'Completed'),
131 (3, 1, 'Bug fixing phase 1', 25.0, '2023-04-01', 'Completed'),
132 (17, 2, 'Grant proposal writing', 35.0, '2023-04-05', 'Completed'),
133 (9, 3, 'Ad copy review', 18.0, '2023-07-01', 'Completed'),
134 (13, 4, 'Sales deck preparation', 22.0, '2023-07-20', 'Completed'),
135 (2, 5, 'Technical specification', 33.0, '2023-09-20', 'In Progress'),
136 (4, 1, 'Final testing', 50.0, '2023-05-15', 'Cancelled'), -- Cancelled task
137 (10, 3, 'Competitor analysis', 30.0, '2023-06-01', 'Completed'),
138 (14, 4, 'Follow-up emails', 20.0, '2023-09-05', 'In Progress'),
139 (18, 2, 'Refine experiment design', 25.0, '2023-05-20', 'In Progress'),
140 (20, 1, 'Security audit', 40.0, '2023-05-20', 'Pending'),
141 (25, 5, 'Deployment planning', 15.0, '2023-12-01', 'Pending'),
142 (2, 1, 'Documentation for design', 20.0, '2023-02-25', 'Completed'),
143 (3, 5, 'Frontend module for Epsilon', 60.0, '2023-10-25', 'In Progress'),
144 (9, 3, 'Press release draft', 25.0, '2023-08-10', 'In Progress'),
145 (13, 4, 'CRM data update', 10.0, '2023-09-10', 'Completed'),
146 (17, 6, 'Research paper outline', 40.0, '2024-01-20', 'Pending'),
147 (4, 1, 'Performance optimization', 0.0, '2023-05-01', 'Pending'), -- Task with 0 hours
148 (6, 2, 'Recruitment for Project Beta', 15.0, '2023-03-15', 'Completed'),
149 (7, 4, 'Onboarding new sales members', 18.0, '2023-07-10', 'In Progress'),
150 (11, 3, 'Video ad script', 30.0, '2023-08-15', 'Pending'),
151 (15, 4, 'Quarterly sales report', 12.0, '2023-09-25', 'Pending'),
152 (19, 6, 'Lab maintenance', 8.0, '2024-02-20', 'Pending'),
153 (22, 3, 'Influencer outreach', 22.0, '2023-07-25', 'In Progress'),
154 (23, 4, 'Legal review of contracts', 16.0, '2023-09-28', 'Pending');

```

Listing 1: Dataset for Company Schema

1 Practice Meanings, Values, Relations, and Advantages

Exercise i.1: Overall Company Metrics

Problem: Calculate the total number of employees, the total salary expenditure, the average salary, the minimum salary, and the maximum salary across all employees. What is the main advantage of using aggregate functions for this?

Exercise i.2: Department Employee Listing

Problem: For each department, list the department name, the number of employees in that department, and a comma-separated list of all employee first names in that department, ordered alphabetically by first name. How does `STRING_AGG` help here?

Exercise i.3: Understanding Different COUNTs

Problem: Find the total number of employees, the number of employees with a `bonus_percentage` recorded, and the number of distinct `performance_rating` values. Explain the difference in meaning for each COUNT.

Exercise i.4: Median Salary and Mode Performance Rating

Problem: Calculate the median salary for employees in the 'Engineering' department and the most common (mode) performance rating for the entire company. What is the value of `PERCENTILE_CONT` and `MODE`?

Exercise i.5: Project Task Hours Distribution

Problem: For each project, display its name, total hours spent, and the variance and standard deviation of hours spent on its tasks. How do `VARIANCE` and `STDDEV` help understand data distribution?

Exercise i.6: Departmental and Cumulative Salaries

Problem: Show the total salary for each department. Also, show the cumulative salary within each department as employees are ordered by their hire date (earliest first). What is the advantage of the window aggregate here?

2 Practice Disadvantages and Potential Pitfalls

Exercise ii.1: Loss of Detail with Average

Problem: Calculate the average salary for the 'Engineering' department. What specific salary information is lost when you only look at this average?

Exercise ii.2: Misleading Aggregate without GROUP BY

Problem: Consider the query `SELECT department_id, MAX(salary) FROM employees;`. Why might this query be misleading or incorrect if the user intends to find the maximum salary *for each department*? What is the potential pitfall?

Exercise ii.3: NULL Handling in AVG()

Problem: Calculate the average `bonus_percentage` for all employees. How does `AVG()` handle NULL values in `bonus_percentage`, and how could this be misleading if not understood?

Exercise ii.4: Aggregate in WHERE Clause

Problem: A manager wants to find departments where the average employee performance rating is below 3.5. They try to write: `SELECT department_id, AVG(performance_rating) FROM employees WHERE AVG(performance_rating) < 3.5 GROUP BY department_id;`. Why will this query fail, and what is the disadvantage or common mistake illustrated here regarding aggregate function placement?

3 Inefficient vs. Efficient Aggregate Usage

Exercise iii.1: Counting Tasks Inefficiently

Problem: A junior analyst needs to find the total number of tasks for 'Project Alpha'. They write a script that fetches all task IDs for 'Project Alpha' into an application and then counts them using application code. Provide the efficient SQL aggregate function solution.

Exercise iii.2: Calculating Average Salary Inefficiently

Problem: To find the average salary of employees hired in 2020, a developer first queries for all salaries of employees hired in 2020, then sums them up and divides by the count in their programming language. How can this be done efficiently in a single SQL query using aggregate functions?

Exercise iii.3: Finding Max Salary Per Department Inefficiently

Problem: A data scientist wants to get a list of departments and, for each, the maximum salary. They write separate queries for each department: `SELECT MAX(salary) FROM employees WHERE department_id = 1;`, then `SELECT MAX(salary) FROM employees WHERE department_id = 2;`, etc. for all departments. Provide a single, efficient SQL query.

Exercise iii.4: Filtering by Total Hours Inefficiently

Problem: An HR assistant needs to find all employees who have logged more than 150 total hours on tasks. They fetch all tasks for every employee, sum the hours in a spreadsheet, and then filter. How can this be done with an efficient SQL query using aggregates and `HAVING`?

4 Hardcore Problem Combining Concepts

Exercise iv.1: Top Employees by Salary per Department with Aggregates and Ranking

Problem: For each department, identify the top 2 employees by salary. For these employees, show their full name, department name, salary, their salary rank within the department (dense rank), the total salary expenditure for their department, and their salary as a percentage of their department's total salary. Only include departments with at least 3 employees. Order results by department name and then by rank.

Exercise iv.2: Project Metrics, Budget Ranking, and Cumulative Budget

Problem: List all projects. For each project, show its name, budget, total hours spent by all employees on that project, and the average hours spent per task on that project. Additionally, rank projects by their budget (highest first). For projects that started in 2023, also show the running total of budgets for projects started in 2023, ordered by their start date.

Exercise iv.3: Employees Above Department Average Salary with Ranking

Problem: For every employee, display their full name, department name, salary, the average salary of their department, and their salary's rank (using `ROW_NUMBER` for unique ranks) within their department. Then, filter this list to show only employees who earn more than their department's average salary and whose hire date is after '2020-01-01'. Order the final result by department name and then by salary in descending order.