

Set Returning Functions & JSON and Array Functions

Data Transformation and Aggregation: Solutions

May 19, 2025

Contents

1	Set Returning Functions (<code>generate_series</code>, <code>unnest</code>) Solutions	5
1.1	(i) Meaning, values, relations (with previous concepts), advantages	5
1.2	(ii) Disadvantages of all its technical concepts	6
1.3	(iii) Cases where people lose advantages due to inefficient solutions	7
1.4	(iv) Hardcore problem combining previous concepts	9
2	JSON and Array Functions Solutions	11
2.1	(i) Meaning, values, relations (with previous concepts), advantages	11
2.2	(ii) Disadvantages of all its technical concepts	12
2.3	(iii) Cases where people lose advantages due to inefficient solutions	14
2.4	(iv) Hardcore problem combining previous concepts	16

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  -- Dataset for Exercises
2  -- Drop tables if they exist to ensure a clean setup
3  DROP TABLE IF EXISTS EventCalendar CASCADE;
4  DROP TABLE IF EXISTS ProjectAssignments CASCADE;
5  DROP TABLE IF EXISTS Employees CASCADE;
6  DROP TABLE IF EXISTS SystemLogs CASCADE;
7  DROP TABLE IF EXISTS ServiceSubscriptions CASCADE;
8
9  -- Create Tables
10 CREATE TABLE Employees (
11     employeeId INT PRIMARY KEY,
12     employeeName VARCHAR(100) NOT NULL,
13     hireDate DATE NOT NULL,
14     department VARCHAR(50),
15     skills TEXT[], -- For unnest, array_append, array_length
16     performanceReviews JSONB -- For JSON functions, e.g., '[{"year": 2022, "rating": 4,
17     "notes": "Good progress"}, ...]'
18 );
19
20 CREATE TABLE ProjectAssignments (
21     assignmentId SERIAL PRIMARY KEY,
22     projectId INT NOT NULL,
23     projectName VARCHAR(100),
24     employeeId INT REFERENCES Employees(employeeId),
25     role VARCHAR(50),
26     assignmentHours INT,
27     assignmentData JSONB -- e.g., '{ "milestones": [{"name": "Phase 1", "status": "
28     completed"}, {"name": "Phase 2", "status": "pending"}], "budget": 5000.00 }'
29 );
30
31 CREATE TABLE SystemLogs (
32     logId SERIAL PRIMARY KEY,
33     logTimestamp TIMESTAMP NOT NULL,
34     serviceName VARCHAR(50),
35     logLevel VARCHAR(10), -- e.g., INFO, ERROR, WARN
36     logDetails JSONB -- e.g., '{ "clientId": "192.168.1.10", "requestPath": "/api/data",
37     "statusCode": 200, "userContext": {"userId": 101, "sessionId": "xyz"} }'
38 );
39
40 CREATE TABLE ServiceSubscriptions (
41     subscriptionId SERIAL PRIMARY KEY,
42     userId INT,
43     customerName VARCHAR(100),
44     serviceType VARCHAR(50),
45     startDate DATE NOT NULL,
46     endDate DATE, -- Can be NULL for ongoing subscriptions
47     monthlyFee DECIMAL(10,2),
48     features JSONB -- e.g., '{ "storageLimitGB": 50, "prioritySupport": true, "addons":
49     ["backup", "monitoring"] }'
50 );
51
52 CREATE TABLE EventCalendar (
53     eventId SERIAL PRIMARY KEY,
54     eventName VARCHAR(100),
55     eventCategory VARCHAR(50),
56     eventStartDate DATE,
57     eventEndDate DATE,
58     expectedAttendees INT,
59     bookedResources TEXT[] -- e.g., '{"Room A", "Projector X"}'
60 );
61
62 -- Populate Tables
63 INSERT INTO Employees (employeeId, employeeName, hireDate, department, skills,
64     performanceReviews) VALUES
65 (1, 'Alice Wonderland', '2020-01-15', 'Engineering', ARRAY['SQL', 'Python', 'Data
66     Analysis'], '[{"year": 2022, "rating": 4, "notes": "Exceeded expectations in Q3"}],
```

```

{"year": 2023, "rating": 5, "notes": "Top performer"}}'),
62 (2, 'Bob The Builder', '2019-03-01', 'Engineering', ARRAY['Java', 'Spring', '
  Microservices'], '[{"year": 2022, "rating": 3, "notes": "Met expectations"}, {"year
  ": 2023, "rating": 4, "notes": "Improved significantly"}]'),
63 (3, 'Charlie Brown', '2021-07-30', 'Sales', ARRAY['Communication', 'Negotiation', 'CRM'
  ], '[{"year": 2023, "rating": 4, "notes": "Good sales figures"}]'),
64 (4, 'Diana Prince', '2018-05-10', 'HR', ARRAY['Recruitment', 'Employee Relations', '
  Legal Knowledge'], NULL),
65 (5, 'Edward Scissorhands', '2022-11-01', 'Engineering', ARRAY['Python', 'Machine
  Learning'], '[{"year": 2023, "rating": 5, "notes": "Innovative solutions"}]'),
66 (6, 'Fiona Gallagher', '2023-02-15', 'Sales', ARRAY['CRM', 'Presentations'], '[]'), --
  Empty JSON array
67 (7, 'George Jetson', '2017-09-01', 'Management', ARRAY['Leadership', 'Strategy'], '[{"
  year": 2022, "rating": 5, "notes": "Excellent leadership"}, {"year": 2023, "rating":
  4, "notes": "Managed team well through transition"}]');
68
69 INSERT INTO ProjectAssignments (projectId, projectName, employeeId, role,
  assignmentHours, assignmentData) VALUES
70 (101, 'Data Warehouse Migration', 1, 'Lead Data Engineer', 120, '{ "milestones": [{"name
  ": "Schema Design", "status": "completed"}, {"name": "ETL Development", "status": "
  in-progress"}], "budget": 20000.00, "critical": true }'),
71 (101, 'Data Warehouse Migration', 2, 'Backend Developer', 80, '{ "milestones": [{"name":
  "API Integration", "status": "pending"}], "budget": 15000.00, "critical": true }'),
72 (102, 'Mobile App Development', 2, 'Lead Mobile Developer', 150, '{ "milestones": [{"
  name": "UI/UX Design", "status": "completed"}, {"name": "Frontend Dev", "status": "
  completed"}, {"name": "Backend Dev", "status": "in-progress"}], "budget": 50000.00,
  "critical": false }'),
73 (103, 'Sales Platform Upgrade', 3, 'Sales Lead', 100, '{ "milestones": [{"name": "
  Requirement Gathering", "status": "completed"}], "budget": 10000.00, "critical":
  false }'),
74 (101, 'Data Warehouse Migration', 5, 'ML Engineer', 60, '{ "milestones": [{"name": "
  Model Training", "status": "in-progress"}], "budget": 12000.00, "critical": true }')
75 (104, 'HR System Implementation', 4, 'HR Specialist', 90, NULL); -- No JSON data
76
77 INSERT INTO SystemLogs (logTimestamp, serviceName, logLevel, logDetails) VALUES
78 ('2023-10-01 10:00:00', 'AuthService', 'INFO', '{ "message": "User login successful", "
  userId": 1, "clientId": "192.168.0.10" }'),
79 ('2023-10-01 10:05:00', 'OrderService', 'ERROR', '{ "message": "Payment processing
  failed", "orderId": 123, "errorCode": "P5001", "details": {"reason": "Insufficient
  funds"} }'),
80 ('2023-10-01 10:10:00', 'ProductService', 'WARN', '{ "message": "Low stock warning", "
  productId": "XYZ123", "currentStock": 5 }'),
81 ('2023-10-02 11:00:00', 'AuthService', 'INFO', '{ "message": "User login successful", "
  userId": 2, "clientId": "192.168.0.15" }'),
82 ('2023-10-02 11:15:00', 'OrderService', 'INFO', '{ "message": "Order placed", "orderId":
  124, "items": ["itemA", "itemB"], "totalAmount": 75.50 }'),
83 (NOW() - INTERVAL '1 day', 'ReportingService', 'DEBUG', '{ "queryId": "q123", "
  executionTimeMs": 1500, "parameters": {"startDate": "2023-01-01", "endDate":
  "2023-01-31"} }');
84
85 INSERT INTO ServiceSubscriptions (userId, customerName, serviceType, startDate, endDate,
  monthlyFee, features) VALUES
86 (101, 'Customer Alpha', 'Premium Cloud Storage', '2023-01-01', '2023-12-31', 20.00, '{ "
  storageLimitGB": 100, "prioritySupport": true, "addons": ["versioning", "encryption
  "] }'),
87 (102, 'Customer Beta', 'Basic VPN Service', '2023-03-15', NULL, 5.00, '{ "dataCapMB":
  5000, "prioritySupport": false, "serverLocations": ["US", "EU"] }'),
88 (103, 'Customer Gamma', 'Standard Streaming', '2023-05-01', '2024-04-30', 10.00, '{ "
  resolution": "1080p", "profiles": 4, "offlineDownload": true }'),
89 (101, 'Customer Alpha', 'Analytics Suite', '2023-06-01', NULL, 50.00, '{ "users": 5, "
  dashboards": 10, "dataSources": ["db1", "s3"] }'),
90 (104, 'Customer Delta', 'Premium Cloud Storage', '2022-11-01', '2023-11-01', 18.00, '{ "
  storageLimitGB": 100, "prioritySupport": true, "addons": ["backup"] }');
91
92 INSERT INTO EventCalendar (eventName, eventCategory, eventStartDate, eventEndDate,
  expectedAttendees, bookedResources) VALUES
93 ('Tech Conference 2024', 'Conference', '2024-03-10', '2024-03-12', 500, ARRAY['Main Hall
  ', 'Audio System', 'Projectors']),
94 ('Product Launch Q1', 'Marketing', '2024-02-15', '2024-02-15', 100, ARRAY['Meeting Room
  Alpha', 'Catering Service']),
95 ('Team Building Workshop', 'HR', '2024-04-05', '2024-04-05', 30, ARRAY['Outdoor Space',
  'Activity Kits']),

```

```
96 ('Quarterly Review Meeting', 'Management', '2024-01-20', '2024-01-20', 15, ARRAY['Board  
Room']),  
97 ('Holiday Party 2023', 'Social', '2023-12-15', '2023-12-15', 150, NULL); -- No resources  
booked
```

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

1 Set Returning Functions (generate_series, unnest) Solutions

1.1 (i) Meaning, values, relations (with previous concepts), advantages

Solution SRF.1.1: Monthly Active Subscription Report

```
1 WITH MonthSeries AS (  
2     SELECT generate_series(  
3         '2023-01-01'::DATE,  
4         '2023-12-01'::DATE,  
5         '1 month'::INTERVAL  
6     ) AS reportMonth  
7 )  
8 SELECT  
9     ss.serviceType,  
10    ms.reportMonth,  
11    COUNT(ss.subscriptionId) AS activeSubscriptions  
12 FROM  
13     MonthSeries ms  
14 CROSS JOIN -- Consider all service types for each month initially  
15     (SELECT DISTINCT serviceType FROM ServiceSubscriptions)  
16     distinctServiceTypes  
17 LEFT JOIN  
18     ServiceSubscriptions ss ON ss.serviceType = distinctServiceTypes.  
19     serviceType  
20     AND ms.reportMonth >= DATE_TRUNC('month', ss.startDate)  
21     AND (ms.reportMonth <= DATE_TRUNC('month', ss.endDate) OR ss.  
22     endDate IS NULL)  
23 GROUP BY  
24     ss.serviceType,  
25     ms.reportMonth  
26 HAVING  
27     ss.serviceType IS NOT NULL -- Only show service types that had  
28     activity  
29 ORDER BY  
30     ss.serviceType,  
31     ms.reportMonth;
```

Solution SRF.1.2: Employee Skills Breakdown

```
1 SELECT  
2     e.employeeName,  
3     e.department,  
4     unnest(e.skills) AS skill  
5 FROM  
6     Employees e  
7 WHERE  
8     e.skills IS NOT NULL AND array_length(e.skills, 1) > 0;
```

1.2 (ii) Disadvantages of all its technical concepts

Solution SRF.2.1: Potential Performance Issue with generate_series

```
1  -- Query that would generate the series (limited for safety)
2  SELECT generate_series(
3      '2023-01-01 00:00:00'::TIMESTAMP,
4      '2023-01-01 00:00:09'::TIMESTAMP, -- Simulating a small part of a
      year
5      '1 second'::INTERVAL
6  ) AS sensorTimestamp
7  LIMIT 10;
8
9  -- Explanation of disadvantage:
10 -- A full year has 365 * 24 * 60 * 60 = 31,536,000 seconds.
11 -- Generating such a large series directly in memory or as an
    intermediate result set
12 -- can consume significant server resources (CPU and RAM) and lead to
    slow query performance
13 -- or even out-of-memory errors. If this data is to be stored, it would
    create a very large table.
14 -- This demonstrates that while powerful, generate_series must be used
    with an understanding
15 -- of the scale of data it can produce.
```

Explanation of disadvantage (as provided in the comment within the code block): A full year has $365 * 24 * 60 * 60 = 31,536,000$ seconds. Generating such a large series directly in memory or as an intermediate result set can consume significant server resources (CPU and RAM) and lead to slow query performance or even out-of-memory errors. If this data is to be stored, it would create a very large table. This demonstrates that while powerful, `generate_series` must be used with an understanding of the scale of data it can produce.

Solution SRF.2.2: Row Explosion with unnest

```
1  -- Hypothetical scenario explanation (no direct query needed to
    demonstrate explosion,
2  -- but we can illustrate a join that would cause it if data were larger
    ):
3  SELECT
4      e.employeeName,
5      s.skill,
6      pa.projectName
7  FROM
8      Employees e,
9      unnest(e.skills) AS s(skill) -- Generates N rows for N skills
10 LEFT JOIN
11     ProjectAssignments pa ON e.employeeId = pa.employeeId -- Could
    multiply by M projects
12 WHERE e.employeeId = 1; -- Example for one employee
13
14 -- Explanation of disadvantage:
15 -- If an employee has 100 skills, 'unnest(e.skills)' for that employee
    creates 100 rows.
16 -- If this employee is also assigned to 10 projects, a subsequent JOIN
    on 'employeeId'
```

```

17 -- (like with 'ProjectAssignments') would result in 100 (skills) * 10 (
    projects) = 1000 rows
18 -- for this single employee *before* any filtering or aggregation
    specific to skill/project combinations.
19 -- This "row explosion" can significantly increase the size of
    intermediate datasets,
20 -- leading to higher memory consumption, slower processing for joins,
    aggregations, and sorting,
21 -- thus degrading overall query performance. It's crucial when using '
    unnest' in conjunction
22 -- with joins to be aware of this multiplicative effect and to apply
    filters as early as possible.

```

Explanation of disadvantage (as provided in the comment within the code block): If an employee has 100 skills, `unnest(e.skills)` for that employee creates 100 rows. If this employee is also assigned to 10 projects, a subsequent JOIN on `employeeId` (like with `ProjectAssignments`) would result in $100 \text{ (skills)} * 10 \text{ (projects)} = 1000$ rows for this single employee *before* any filtering or aggregation specific to skill/project combinations. This "row explosion" can significantly increase the size of intermediate datasets, leading to higher memory consumption, slower processing for joins, aggregations, and sorting, thus degrading overall query performance. It's crucial when using `unnest` in conjunction with joins to be aware of this multiplicative effect and to apply filters as early as possible.

1.3 (iii) Cases where people lose advantages due to inefficient solutions

Solution SRF.3.1: Inefficiently Generating Date Sequences

```

1 WITH RECURSIVE DateSequenceCTE AS (
2     SELECT '2024-01-01'::DATE AS aDate
3     UNION ALL
4     SELECT (aDate + INTERVAL '1 day')::DATE
5     FROM DateSequenceCTE
6     WHERE aDate < '2024-01-31'::DATE
7 )
8 SELECT aDate FROM DateSequenceCTE;

```

```

1 SELECT generate_series(
2     '2024-01-01'::DATE,
3     '2024-01-31'::DATE,
4     '1 day'::INTERVAL
5 ) AS aDate;

```

Advantage Explanation: The `generate_series` approach is more advantageous because:

1. **Conciseness:** It's a single function call, making the SQL much shorter and easier to read and understand.
2. **Performance:** `generate_series` is typically implemented as a highly optimized internal function, often outperforming recursive CTEs for simple sequence generation.

3. **Simplicity:** It directly expresses the intent of generating a series without the boilerplate of setting up a recursive query (anchor member, recursive member, termination condition).

People might use recursive CTEs if they are unaware of `generate_series` or if they need more complex, conditional series generation that `generate_series` alone can't handle (though often, `generate_series` can be combined with other logic).

Solution SRF.3.2: Inefficiently Handling Array Elements

```
1 SELECT
2     employeeName ,
3     skills
4 FROM
5     Employees
6 WHERE
7     employeeId = 1 AND
8     array_to_string(skills, ',') LIKE '%Python%'; -- Prone to errors (e
    .g. skill 'Pythonista') and inefficient
```

```
1 -- Using array contains operator @>
2 SELECT
3     employeeName ,
4     skills
5 FROM
6     Employees
7 WHERE
8     employeeId = 1 AND
9     skills @> ARRAY['Python'];
10
11 -- Alternative using ANY with unnest (conceptually, though often
    optimized without full unnesting)
12 SELECT
13     e.employeeName ,
14     e.skills
15 FROM
16     Employees e
17 WHERE
18     e.employeeId = 1 AND
19     'Python' = ANY(e.skills);
20 -- Or with explicit unnest for checking:
21 SELECT DISTINCT e.employeeName, e.skills
22 FROM Employees e, unnest(e.skills) s
23 WHERE e.employeeId = 1 AND s = 'Python';
```

Advantage Explanation: The SQL-native array approaches (`@>`, `ANY`, or `unnest` with `WHERE`) are advantageous because:

1. **Accuracy:** They correctly match whole array elements, avoiding partial matches or issues with delimiters that `array_to_string` and `LIKE` might suffer from (e.g., searching for `'Py'` would match `'Python'`).

2. **Performance:** PostgreSQL has optimized operators and indexing capabilities (GIN indexes) for arrays, making these searches much faster than string operations on concatenated array content, especially for large arrays or tables.
3. **Readability and Intent:** Using array-specific operators clearly conveys the intent to operate on array elements.

People might resort to string manipulation if they are unfamiliar with array functions and operators, leading to less robust and slower queries.

1.4 (iv) Hardcore problem combining previous concepts

Solution SRF.4.1: Comprehensive Project Health and Skill Utilization Report

```

1 WITH MonthlySeries AS (
2     SELECT generate_series(
3         '2023-01-01'::DATE,
4         '2023-06-01'::DATE,
5         '1 month'::INTERVAL
6     ) AS reportMonth
7 ),
8 ProjectEmployeeSkills AS (
9     SELECT
10         pa.projectId,
11         pa.projectName,
12         e.employeeId,
13         e.hireDate,
14         s.skill,
15         pa.assignmentHours,
16         (pa.assignmentData ->> 'critical')::BOOLEAN AS isCritical
17     FROM
18         ProjectAssignments pa
19     JOIN
20         Employees e ON pa.employeeId = e.employeeId
21     CROSS JOIN
22         unnest(e.skills) AS s(skill) -- unnest skills
23 ),
24 ProjectsWithPythonSkill AS ( -- Identify projects that utilize 'Python'
25     skill
26     SELECT DISTINCT projectId
27     FROM ProjectEmployeeSkills
28     WHERE skill = 'Python'
29 )
30 SELECT
31     ms.reportMonth,
32     pes.projectName,
33     SUM(pes.assignmentHours) AS totalAssignedHours,
34     MAX(CASE WHEN pes.isCritical THEN 1 ELSE 0 END)::BOOLEAN AS
criticalProjectFlag, -- True if any assignment is critical
35     STRING_AGG(DISTINCT pes.skill, ',' ORDER BY pes.skill) AS
listOfDistinctSkillsUtilized,
36     ROUND(AVG(
37         EXTRACT(YEAR FROM AGE(ms.reportMonth, pes.hireDate)) +
        EXTRACT(MONTH FROM AGE(ms.reportMonth, pes.hireDate)) / 12.0 +

```

```

38         EXTRACT(DAY FROM AGE(ms.reportMonth, pes.hireDate)) / 365.25 --
39         More precise age calculation
40     ), 2) AS averageYearsOfService
41 FROM
42     MonthlySeries ms
43     CROSS JOIN -- To consider each project for each month
44     (SELECT DISTINCT projectId, projectName FROM ProjectAssignments)
45     uniqueProjects
46 JOIN
47     ProjectEmployeeSkills pes ON uniqueProjects.projectId = pes.
48     projectId
49 WHERE
50     pes.hireDate <= ms.reportMonth -- Employee active in the report
51     month
52     AND pes.projectId IN (SELECT projectId FROM ProjectsWithPythonSkill
53     ) -- Filter for projects with Python skill
54 GROUP BY
55     ms.reportMonth,
56     pes.projectName,
57     pes.projectId -- Added projectId to GROUP BY for correctness if
58     projectName is not unique across IDs
59 ORDER BY
60     ms.reportMonth,
61     pes.projectName;

```

2 JSON and Array Functions Solutions

2.1 (i) Meaning, values, relations (with previous concepts), advantages

Solution JAF.1.1: Extracting Specific Log Information

```
1 SELECT
2     logTimestamp,
3     logDetails ->> 'clientId' AS clientId,
4     logDetails -> 'userContext' ->> 'userId' AS userId
5     -- Alternative for userId: jsonb_extract_path_text(logDetails, '
6     userContext', 'userId')
7 FROM
8     SystemLogs
9 WHERE
10    serviceName = 'AuthService' AND logLevel = 'INFO';
```

Solution JAF.1.2: Expanding Performance Review Details

```
1 SELECT
2     e.employeeName,
3     (review.reviewData ->> 'year')::INT AS reviewYear,
4     (review.reviewData ->> 'rating')::INT AS reviewRating,
5     review.reviewData ->> 'notes' AS reviewNotes
6 FROM
7     Employees e,
8     jsonb_array_elements(e.performanceReviews) AS review(reviewData) --
9     reviewData is each JSON object in the array
10 WHERE
11     e.performanceReviews IS NOT NULL AND jsonb_typeof(e.
12     performanceReviews) = 'array' AND jsonb_array_length(e.
13     performanceReviews) > 0;
```

Solution JAF.1.3: Constructing a Simplified Project Overview JSON

```
1 SELECT
2     pa.projectId,
3     pa.projectName, -- Keep projectName separate or include in
4     jsonb_build_object
5     (
6         'projectNameForJson', pa.projectName,
7         'assignedEmployeeIds', jsonb_agg(DISTINCT pa.employeeId ORDER
8         BY pa.employeeId) -- Using jsonb_agg for a JSON array
9     ) AS projectOverviewJson
10 FROM
11     ProjectAssignments pa
12 GROUP BY
13     pa.projectId, pa.projectName
14 ORDER BY
15     pa.projectId;
16 -- If a SQL array within the JSON is desired:
17 -- SELECT
18 --     pa.projectId,
```

```

18 --      jsonb_build_object(
19 --          'projectNameForJson', pa.projectName,
20 --          'assignedEmployeeIds', TO_JSONB(ARRAY_AGG(DISTINCT pa.
21 --              employeeId ORDER BY pa.employeeId))
22 --      ) AS projectOverviewJson
23 -- FROM
24 --     ProjectAssignments pa
25 -- GROUP BY
26 --     pa.projectId, pa.projectName
27 -- ORDER BY
28 --     pa.projectId;

```

Solution JAF.1.4: Updating Event Resources and Checking Count

```

1 -- This exercise implies an update then select, or a select simulating
2 -- the update.
3 -- For an exercise, simulating the update in a SELECT is safer.
4 -- If an actual UPDATE is desired:
5 -- UPDATE EventCalendar
6 -- SET bookedResources = array_append(bookedResources, 'WiFi Access
7 -- Point')
8 -- WHERE eventName = 'Tech Conference 2024';
9
10 SELECT
11     eventName,
12     array_append(COALESCE(bookedResources, '{}'), 'WiFi Access Point')
13 AS updatedResources,
14     array_length(array_append(COALESCE(bookedResources, '{}'), 'WiFi
15     Access Point'), 1) AS totalResources
16 FROM
17     EventCalendar
18 WHERE
19     eventName = 'Tech Conference 2024';

```

2.2 (ii) Disadvantages of all its technical concepts

Solution JAF.2.1: Performance of Complex JSON Queries vs. Normalized Data

Explanation: While JSONB is efficient for storage and flexible, querying deeply nested values (e.g., logDetails -> 'path1' -> 'path2' -> ... -> 'targetValue') has disadvantages compared to dedicated relational columns:

1. **Query Performance:** Accessing deeply nested elements requires traversing the JSON structure. Even with GIN/GiST indexes on the JSONB column, which significantly help, direct access to a value in its own indexed column is typically faster, especially for highly selective queries on that specific piece of data. The overhead of JSON parsing/traversal, though optimized, can be higher than a direct column read.
2. **Indexing Complexity:** While GIN/GiST indexes can index JSONB paths or all keys, designing and maintaining these indexes for optimal performance across various deep-path queries can be more complex than indexing simple scalar columns. Full JSONB indexing can also lead to larger index sizes.

3. **Readability and Maintenance:** Queries involving multiple chained `->` or `->>` operators for deep paths can become less readable and harder to maintain than queries on flat table structures.
4. **Data Integrity and Typing:** Relational columns enforce data types at the column level. While JSONB preserves types, extracting and casting them correctly from deep within JSON requires careful query construction and offers less schema enforcement than separate columns.

If certain deeply nested values are frequently queried and critical for performance, it might be more advantageous to extract them into their own relational columns, even if it means some data duplication or a hybrid model.

Solution JAF.2.2: Array Overuse and Normalization

Explanation: Storing complex skill information (skill name, level, experience) as encoded strings within a single `TEXT[]` array (e.g., `ARRAY['Python:Expert:5yrs', 'SQL:Intermediate:3yrs']`) has several disadvantages compared to a normalized approach (e.g., a separate `EmployeeSkills` table with columns `employeeId`, `skillName`, `skillLevel`, `yearsExperience`):

1. **Querying Complexity:** To query for employees with 'Expert' level in 'Python', you'd need complex string parsing (e.g., `SUBSTRING`, `SPLIT_PART`) within the array elements, which is inefficient and error-prone. A separate table allows simple `WHERE skillName = 'Python' AND skillLevel = 'Expert'`.
2. **Data Integrity and Consistency:** There's no easy way to enforce that the skill level is always one of 'Beginner', 'Intermediate', 'Expert', or that years of experience is a number. Encoded strings are free-form. A separate table can have foreign keys, check constraints, and appropriate data types.
3. **Updating Difficulty:** Updating just the skill level for one skill for an employee becomes a complex array manipulation task (find element, parse, change part, reconstruct string, update array). In a separate table, it's a simple `UPDATE` on a specific row.
4. **Indexing Limitations:** Indexing attributes within these encoded strings is very difficult. A separate table allows indexing on `skillName`, `skillLevel`, etc., for fast lookups.
5. **Violation of First Normal Form (1NF):** While the array itself might contain atomic text values, if those text values are composite (like `'Python:Expert:5yrs'`), it's a way of smuggling multiple logical attributes into a single array element. True 1NF would suggest each attribute (skill name, level, experience) be in its own column, often necessitating a separate table for a one-to-many or many-to-many relationship.

Using arrays for simple lists of atomic values is fine (like current `skills` which are just names), but for structured data within each "element," a normalized approach is generally superior for data integrity, querying, and maintenance.

2.3 (iii) Cases where people lose advantages due to inefficient solutions

Solution JAF.3.1: Inefficiently Querying JSON Data with String Matching

```
1 SELECT
2     logId,
3     logDetails
4 FROM
5     SystemLogs
6 WHERE
7     logDetails::TEXT LIKE '%"orderId": 123%'; -- Inefficient and error-
        prone (e.g. whitespace, "orderId": "123")
```

```
1 SELECT
2     logId,
3     logDetails
4 FROM
5     SystemLogs
6 WHERE
7     logDetails ->> 'orderId' = '123'; -- Assumes orderId value is
        stored as a number or string '123'
8     -- If orderId is always a number in JSONB:
9     -- WHERE (logDetails -> 'orderId')::TEXT = '123'; -- JSONB numbers
        need casting to text for = with text
10    -- Or, the most robust way for JSONB if orderId is a number:
11    -- WHERE logDetails @> '{"orderId": 123}'::JSONB;
```

Advantage Explanation: The JSONB-specific approach (`->>`, `@>`) is advantageous because:

1. **Accuracy and Type Safety:** JSONB operators understand JSON structure and types. `logDetails ->> 'orderId' = '123'` correctly extracts the value associated with the key 'orderId' and compares it as text. `logDetails @> '{"orderId": 123}'::JSONB` correctly checks for the key and number value. `LIKE '%"orderId": 123%'` is a fragile string match; it could falsely match if 123 appeared as part of another value or if there was different whitespace (e.g., `"orderId" : 123`). It also doesn't distinguish between the number 123 and the string "123" in the JSON if not careful with the `LIKE` pattern.
2. **Performance:** PostgreSQL's JSONB operators are highly optimized and can leverage GIN/GiST indexes on JSONB columns. Casting to `TEXT` and using `LIKE` prevents effective use of JSONB-specific indexes and involves costly string operations on potentially large JSON structures.
3. **Readability and Intent:** Using JSONB operators clearly communicates that the query is interacting with JSON data structure, not just arbitrary text.

Solution JAF.3.2: Storing Multiple Flags as a Comma-Separated String Instead of JSON/Array

```

1 -- Hypothetical table structure: CREATE TABLE Users (userId INT,
   userTags VARCHAR(255));
2 -- INSERT INTO Users VALUES (1, 'active,premium,verified'), (2, 'active
   ,verified'), (3, 'premium');
3
4 -- Assuming 'Users' table exists and is populated as above:
5 -- SELECT userId, userTags
6 -- FROM Users -- Hypothetical table
7 -- WHERE
8 --     (',' || userTags || ',') LIKE '%,premium,%'; -- Common workaround
   for whole-word matching

```

Advantageous Approaches (Array or JSONB):

1. Using TEXT[] (Array): If userTags was TEXT[] (e.g., ARRAY['active', 'premium', 'verified']):

```

1 -- Query for 'premium' tag:
2 -- SELECT userId, userTags FROM Users WHERE userTags @> ARRAY['premium
   '];
3 -- Or: SELECT userId, userTags FROM Users WHERE 'premium' = ANY(
   userTags);

```

Advantages:

- **Atomicity:** Each tag is a distinct element.
- **Efficient Searching:** Array operators (@>, &&, ANY) are optimized and can use GIN indexes.
- **Easier Manipulation:** Adding/removing tags is cleaner with array_append, array_remove.

2. Using JSONB (Object for flags): If userTags was JSONB (e.g., {"active": true, "premium": true, "verified": true} or {"tags": ["active", "premium", "verified"]}):

```

1 -- Query for 'premium' flag (if keys are flags):
2 -- SELECT userId, userTags FROM Users WHERE (userTags ->> 'premium')::
   BOOLEAN = TRUE;
3 -- Or using existence operator if premium key means it's set:
4 -- SELECT userId, userTags FROM Users WHERE userTags ? 'premium';
5
6 -- If JSONB stores an array: {"tags": ["active", "premium", "verified
   "]}
7 -- SELECT userId, userTags FROM Users WHERE userTags -> 'tags' @> '['"
   premium"'::JSONB;

```

Advantages:

- **Structured Data:** Can store key-value pairs (e.g., flag name and its boolean state) or arrays within JSON.
- **Powerful Querying:** JSON operators and path expressions offer flexible querying, also benefiting from GIN indexes.
- **Schema Flexibility:** Easily add new flags/tags without altering table structure.

Overall Advantages of Array/JSONB over Comma-Separated Strings:

- **Data Integrity:** Avoids issues with delimiters within tag names or inconsistent formatting.
- **Search Performance:** SQL-native array/JSONB operations are significantly faster and can be indexed effectively. String LIKE on concatenated strings is slow and cannot use specialized indexes well for this purpose.
- **Maintainability:** Easier to add, remove, or modify individual tags/flags without complex string parsing and reconstruction.
- **Clarity of Intent:** Using appropriate data types makes the data model and queries more understandable.

People often use comma-separated strings due to perceived simplicity initially or when migrating from systems that use such formats, but they lose out on the robustness, performance, and maintainability offered by PostgreSQL's array and JSON types.

2.4 (iv) Hardcore problem combining previous concepts

Solution JAF.4.1: Advanced Customer Subscription Feature Analysis and Aggregation

```

1 WITH CustomerSubscriptionDetails AS (
2     SELECT
3         ss.subscriptionId,
4         ss.customerName,
5         ss.serviceType,
6         ss.monthlyFee,
7         ss.startDate,
8         ss.endDate,
9         ss.features,
10        CASE
11            WHEN ss.endDate IS NULL OR ss.endDate > CURRENT_DATE THEN '
Active'
12            ELSE 'Expired'
13        END AS status,
14        EXTRACT(YEAR FROM AGE(COALESCE(ss.endDate, CURRENT_DATE), ss.
startDate)) * 12 +
15        EXTRACT(MONTH FROM AGE(COALESCE(ss.endDate, CURRENT_DATE), ss.
startDate)) AS durationMonths,
16        COALESCE((ss.features ->> 'prioritySupport')::BOOLEAN, FALSE)
AS hasPrioritySupport,
17        -- Extract addons array; ensure it's an array and handle NULLs/
non-arrays gracefully
18        CASE
19            WHEN jsonb_typeof(ss.features -> 'addons') = 'array'
20            THEN (SELECT jsonb_agg(elem) FROM jsonb_array_elements_text
(ss.features -> 'addons') AS elem)
21            ELSE '[]'::JSONB
22        END AS featureListJsonArray
23    FROM
24        ServiceSubscriptions ss
25 ),
26 CustomersWithPrioritySupport AS (
27     SELECT DISTINCT customerName

```



```

28     FROM CustomerSubscriptionDetails
29     WHERE hasPrioritySupport = TRUE
30 )
31 SELECT
32     csd.customerName, -- Included for clarity, though it's in the JSON
33     jsonb_build_object(
34         'customerName', csd.customerName,
35         'totalMonthlyFee', COALESCE(SUM(csd.monthlyFee) FILTER (WHERE
36             csd.status = 'Active'), 0.00),
37         'activeServicesCount', COUNT(*) FILTER (WHERE csd.status = '
38             Active'),
39         'serviceDetails', jsonb_agg(
40             jsonb_build_object(
41                 'serviceType', csd.serviceType,
42                 'status', csd.status,
43                 'durationMonths', csd.durationMonths,
44                 'featureList', csd.featureListJsonArray,
45                 'hasPrioritySupport', csd.hasPrioritySupport
46             ) ORDER BY csd.startDate DESC
47         )
48     ) AS customerReport
49 FROM
50     CustomerSubscriptionDetails csd
51 WHERE
52     csd.customerName IN (SELECT customerName FROM
53         CustomersWithPrioritySupport)
54 GROUP BY
55     csd.customerName
56 ORDER BY
57     SUM(csd.monthlyFee) FILTER (WHERE csd.status = 'Active') DESC NULLS
58     LAST;

```