# Practical SQL

# Referring to SQL results in SQL queries -- Subqueries

# Subqueries embed one query within another

## These embedded queries are 'subqueries'

We have to watch the table returned by the subquery carefully

The returned table has to have the fields and row expected of it by the calling query

Two varieties, 'correlated' and 'uncorrelated'

'correlated' is cooler but harder, we'll do 'uncorrelated' first

# Where we left off:

## Say we want a list of those presidents from the three states that have sent the most presidents

We can get the list by a query:

```
 1  SELECT
 2      state,
 3      COUNT(*) AS state_count
 4  FROM
 5      sampdb.president
 6  GROUP BY
 7      state
 8  ORDER BY
 9      state_count DESC
10  LIMIT 3
11  ;
12
```

```
 1  +-------+-------------+
 2  | state | state_count |
 3  +-------+-------------+
 4  | VA    |           8 |
 5  | OH    |           7 |
 6  | MA    |           4 |
 7  +-------+-------------+
 8
```

# First, an easier example

## Say we want a list of students scoring above on grade event 3

### Getting the score is easy:

```
 1  SELECT
 2      event_id,
 3      AVG(score)
 4  FROM
 5      sampdb.score
 6  WHERE
 7          event_id = 3
 8  GROUP BY
 9      event_id
10  LIMIT 5
11  ;
12
```

### Filtering is easy

```
 1  SELECT
 2      student_id AS student,
 3      event_id AS event,
 4      score
 5  FROM
 6      sampdb.score
 7  WHERE
 8          event_id = 3
 9      AND score > 78.2258
10  ;
11
```

### Results

```
 1  +---------+-------+-------+
 2  | student | event | score |
 3  +---------+-------+-------+
 4  |       1 |     3 |    88 |
 5  |       2 |     3 |    84 |
 6  |       5 |     3 |    97 |
 7  |       6 |     3 |    83 |
 8  |       7 |     3 |    88 |
 9  +---------+-------+-------+
10
```

### But:

What if the data change?

We want results for another event?

# Better: refer directly to the average

```
 1  SELECT
 2      student_id,
 3      event_id,
 4      score
 5  FROM
 6      sampdb.score
 7  WHERE
 8          event_id = 3
 9      AND score > (
10          SELECT
11              AVG(score)
12          FROM
13              sampdb.score
14          WHERE
15              event_id = 3
16          GROUP BY
17              event_id
18      )
19  LIMIT 5;
20
```

```
 1  +------------+-------+
 2  | student_id | score |
 3  +------------+-------+
 4  |          1 |    88 |
 5  |          2 |    84 |
 6  |          5 |    97 |
 7  |          6 |    83 |
 8  |          7 |    88 |
 9  +------------+-------+
10
```

**Now:**

- syntax is more complex

- but query automatically changes output for any data change

- we can get a similar list simply by changing the event_id we are seeking

- BUT -- what if we change event_id in one place, and not the other . . .?

# Better still: Remove duplication

```
 1  SELECT
 2    student_id,
 3    score
 4  FROM
 5    sampdb.score AS scr
 6  WHERE
 7    event_id = 3
 8    AND score > (
 9      SELECT
10        AVG(score)
11      FROM
12        sampdb.score
13      WHERE
14        event_id=scr.event_id
15      GROUP BY
16        event_id
17    )
18  LIMIT 5;
19
```

```
 1  +------------+-------+
 2  | student_id | score |
 3  +------------+-------+
 4  |          1 |    88 |
 5  |          2 |    84 |
 6  |          5 |    97 |
 7  |          6 |    83 |
 8  |          7 |    88 |
 9  +------------+-------+
10
```

**Now:**

- The subquery's event_id filter looks to each particular record for the value of event_id used in its filter
- Thus our query only refers to the desired event_id once
- e.g., we cannot make the 'event_id values out of step' mistake
- This is a 'correlated subquery', and we will spend more time on them later

# Back to the presidents problem

## Revise our condition for state inclusion

This will automatically update values for changes in the data, or application to another data set

It is also easier to change the logic -- if we want the top 5 states, we change one value and the implications flow
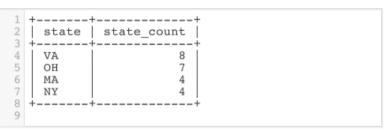
```
 1  SELECT
 2    MIN(st_count)
 3  FROM
 4  (
 5    SELECT
 6      state,
 7      COUNT(*) AS st_count
 8    FROM
 9      sampdb.president
10    GROUP BY
11      state
12    ORDER BY
13      state_count DESC
14    LIMIT 3
15  ) AS tmp
16  ;
17
```

```
 1  +-----------------+
 2  | MIN(st_count) |
 3  +-----------------+
 4  |               4 |
 5  +-----------------+
 6
```

# Revise listing query

## Uses presidents count criteria for state inclusion

```
 1 SELECT
 2   state, COUNT(*) as state_count
 3 FROM
 4   sampdb.president
 5 GROUP BY
 6   state
 7 HAVING
 8   COUNT(*) >=
 9   (
10     SELECT
11       MIN(state_count)
12     FROM
13       (
14         SELECT
15           state,
16           COUNT(*) AS state_count
17         FROM
18           sampdb.president
19         GROUP BY
20           state
21         ORDER BY
22           state_count DESC
23         LIMIT 3
24       ) AS t3
25   )
26 ORDER BY state_count DESC
27 ;
28
```

```
1 +-------+-------------+
2 | state | state_count |
3 +-------+-------------+
4 |  VA   |           8 |
5 |  OH   |           7 |
6 |  MA   |           4 |
7 |  NY   |           4 |
8 +-------+-------------+
9
```

- This gives us a 'dynamic' derivation of the state list we need

- HAVING

- - similar to WHERE

- - appears after GROUP BY

- - refers to, and filters on, the grouped records and aggregated fields

- So WHERE filters the records that are aggregated into groups, and HAVING filters the groups themselves

# Putting all this together provides our answer

```
 1  SELECT last_name, state
 2  FROM sampdb.president
 3  WHERE
 4    state IN
 5    (
 6      SELECT state
 7      FROM sampdb.president
 8      GROUP BY state
 9      HAVING
10        COUNT(*) >=
11        (
12          SELECT
13            MIN(state_count)
14          FROM
15            (
16              SELECT COUNT(*) AS state_count
17              FROM sampdb.president
18              GROUP BY state
19              ORDER BY state_count DESC
20              LIMIT 3
21            ) AS t3
22        )
23    )
24  ORDER BY
25    state, last_name, first_name
26  ;
27
```

```
 1  +------------+-------+
 2  | last_name  | state |
 3  +------------+-------+
 4  | Adams      | MA    |
 5  | Adams      | MA    |
 6  | Bush       | MA    |
 7  | Kennedy    | MA    |
 8  | Fillmore   | NY    |
 9  | Roosevelt  | NY    |
10  | Roosevelt  | NY    |
11  | Van Buren  | NY    |
12  | Garfield   | OH    |
13  | Grant      | OH    |
14  | Harding    | OH    |
15  | Harrison   | OH    |
16  | Hayes      | OH    |
17  | McKinley   | OH    |
18  | Taft       | OH    |
19  | Harrison   | VA    |
20  | Jefferson  | VA    |
21  | Madison    | VA    |
22  | Monroe     | VA    |
23  | Taylor     | VA    |
24  | Tyler      | VA    |
25  | Washington | VA    |
26  | Wilson     | VA    |
27  +------------+-------+
28
```

- Inner query lists the top three states and their counts
- Next extracts the minimum value of those counts
- Next out lists the states with that count or higher
- Last query takes presidents whose states fall in that list

# Subqueries in SELECT field expressions

```
 1  SELECT
 2  ( SELECT
 3      ROUND(AVG(score),1)
 4    FROM sampdb.score
 5    WHERE event_id = 1
 6    GROUP BY event_id
 7  )   AS 1_scr,
 8  ( SELECT
 9      ROUND(AVG(score),1)
10    FROM sampdb.score
11    WHERE event_id = 2
12    GROUP BY event_id
13  )   AS 2_scr,
14  ( SELECT
15      ROUND(AVG(score),1)
16    FROM sampdb.score
17    WHERE event_id = 3
18    GROUP BY event_id
19  )   AS 3_scr
20
21  -- Notice the SELECT
22  -- expression is
23  -- just a field list
24
25  ;
26
27
28
```

```
+-------+-------+-------+
| 1_scr | 2_scr | 3_scr |
+-------+-------+-------+
|  15.1 |  14.2 |  78.2 |
+-------+-------+-------+
```

**Subqueries can be used in field expressions**

- Here we use them to replicate an Excel pivot table
- Notice that this doesn't handle "arbitrary" columns
- We have to designate a field in our output for each column
- The fundamental problem is we can specify sets of records, but not fields
- There are hacks around this, but they're complicated

# Field Subqueries used for calculations

```sql
SELECT
    score, student_id AS stdnt, event_id AS vnt,
    event_average AS vnt_avg, std_dev AS std,
    ( score - event_average ) / std_dev AS z_score
FROM
(
    SELECT
        s1.score, s1.student_id, s1.event_id,
        (   SELECT AVG(score)
            FROM sampdb.score AS s2
            WHERE s2.event_id = s1.event_id
            GROUP BY event_id
        ) AS event_average,
        (   SELECT STD(score)
            FROM sampdb.score AS s2
            WHERE s2.event_id = s1.event_id
            GROUP BY event_id
        ) AS std_dev
    FROM sampdb.score AS s1
) AS tmp

LIMIT 5
;
```

```
+-------+-------+-----+---------+--------+-------------+
| score | stdnt | vnt | vnt_avg | std    | z_score     |
+-------+-------+-----+---------+--------+-------------+
|    20 |     1 |   1 | 15.1379 | 3.7849 |  1.28458961 |
|    20 |     3 |   1 | 15.1379 | 3.7849 |  1.28458961 |
|    18 |     4 |   1 | 15.1379 | 3.7849 |  0.75618024 |
|    13 |     5 |   1 | 15.1379 | 3.7849 | -0.56484320 |
|    18 |     6 |   1 | 15.1379 | 3.7849 |  0.75618024 |
+-------+-------+-----+---------+--------+-------------+
```

# Correlated Subqueries work record by record

## The 'inner' queries are run for each record in the 'outer' query

### Definitions

- 'Outer' query is the query calling on result
- 'Inner' query supplies result
- Records in separate are referenced by aliases

### Outer query checks inner query for each of its own records

1. Outer gets record
2. Outer consults inner, 'correlated' query
3. Inner query takes values from record in process for outer query
4. Inner passes result out to outer query

### Are duplicative inner queries cached?

- It depends
- Performance is complex, requires first mastering syntax

# HAVING

## Seen earlier in subquery

### HAVING applies criteria to output of aggregations

- Calculated after WHERE and GROUP BY
- WHERE filters input records to initial query
- HAVING filters resultant aggregations

# HAVING Example

## States with more than one President whose last name doesn't begin with a vowel

```
 1  SELECT
 2    state,
 3    COUNT(*) AS state_count
 4  FROM
 5    sampdb.president
 6  WHERE
 7    LEFT(last_name, 1)
 8    NOT IN
 9    ('A', 'E', 'I', '
10      O', 'U')
11  GROUP BY
12    state
13  HAVING
14    state_count >= 2
15  ;
16
17
```

```
 1  +-------+-------------+
 2  | state | state_count |
 3  +-------+-------------+
 4  | MA    |           2 |
 5  | NC    |           2 |
 6  | NY    |           4 |
 7  | OH    |           7 |
 8  | VA    |           8 |
 9  +-------+-------------+
10
```

- WHERE clause strips out the vowel last names
- HAVING gets the states with more than one (now not-vowel) president
- Splitting up the WHERE and HAVING criteria gives a finer control over the output set

# Queries for Analysis

# Using queries for analysis

- Notice that all of our queries return (e.g., output) tables
  - The output is a set of records, with defined fields
- We can perform further operations on those returned tables, and join them to one another
- By cascading queries, we can filter and massage our data to get where we want
  - "views" give us a tool for holding queries "in place"

# Create a "base table"

```
DROP VIEW IF EXISTS sampdb.president_age;

CREATE VIEW
    sampdb.president_age
AS
    SELECT
        last_name,
        state,
        ROUND(DATEDIFF(death, birth) / 365, 1) as AGE
    FROM
        sampdb.president;
```

```
+------------+-------+------+
| last_name  | state | age  |
+------------+-------+------+
| Washington | VA    | 67.9 |
| Adams      | MA    | 90.7 |
| Jefferson  | VA    | 83.3 |
| Madison    | VA    | 85.3 |
| Monroe     | VA    | 73.2 |
| Adams      | MA    | 80.7 |
| Jackson    | SC    | 78.3 |
| Van Buren  | NY    | 79.7 |
| Harrison   | VA    | 68.2 |
| Tyler      | VA    | 71.9 |
| Polk       | NC    | 53.7 |
| Taylor     | VA    | 65.7 |
| Fillmore   | NY    | 74.2 |
```

- "Stores" the result of the SELECT query in a "table" available to other queries
- Uses functions to create an << age >> column from birth and death

# Derive a second table from our base

```
DROP VIEW IF EXISTS sampdb.president_aggregates;

CREATE VIEW
   sampdb.president_aggregates
AS
   SELECT
      state,
      COUNT(state) as StateCount,
      AVG(age) as AverageAge,
      MAX(age) as MaxAge,
      MIN(age) as MinAge
   FROM
      sampdb.president_age
   GROUP BY
      state
   ORDER BY
      StateCount DESC
;
```

- Uses our prior query to create a table of aggregated values

| state | StateCount | AverageAge | MaxAge | MinAge |
|-------|-----------|------------|--------|--------|
| VA | 8 | 72.83750 | 85.3 | 65.7 |
| OH | 7 | 62.87143 | 72.5 | 49.9 |
| MA | 4 | 72.63333 | 90.7 | 46.5 |
| NY | 4 | 69.32500 | 79.7 | 60.2 |
| TX | 2 | 71.45000 | 78.5 | 64.4 |
| NC | 2 | 60.15000 | 66.6 | 53.7 |

# Combine the two to get our result

```
SELECT
    sub.state,
    sub.last_name,
    sub.age
FROM
    sampdb.president_age sub
INNER JOIN
    sampdb.president_aggregates sub2
ON
    sub.state = sub2.state
    AND sub.age = sub2.MaxAge
```

- Now we can join our Age table with our Aggregates table to get oldest presidents by State

| state | last_name | age |
|-------|-----------|------|
| MA    | Adams     | 90.7 |
| VA    | Madison   | 85.3 |
| SC    | Jackson   | 78.3 |
| NY    | Van Buren | 79.7 |
| NH    | Pierce    | 64.9 |
| PA    | Buchanan  | 77.2 |

# Cascading through queries to result

| last_name | first_name | suffix | city | state | birth | death |
|-----------|------------|--------|------|-------|-------|-------|
| Washington | George | NULL | Wakefield | VA | 1732-02-22 | 1799-12-14 |
| Adams | John | NULL | Braintree | MA | 1735-10-30 | 1826-07-04 |
| Jefferson | Thomas | NULL | Albemarle County | VA | 1743-04-13 | 1826-07-04 |
| Madison | James | NULL | Port Conway | VA | 1751-03-16 | 1836-06-28 |
| Monroe | James | NULL | Westmoreland County | VA | 1758-04-28 | 1831-07-04 |
| Adams | John Quincy | NULL | Braintree | MA | 1767-07-11 | 1848-02-23 |

| last_name | state | age |
|-----------|-------|-----|
| Washington | VA | 67.9 |
| Adams | MA | 90.7 |
| Jefferson | VA | 83.3 |
| Madison | VA | 85.3 |
| Monroe | VA | 73.2 |
| Adams | MA | 80.7 |
| Jackson | SC | 78.3 |
| Van Buren | NY | 79.7 |
| Harrison | VA | 68.2 |
| Tyler | VA | 71.9 |
| Polk | NC | 53.7 |
| Taylor | VA | 65.7 |
| Fillmore | NY | 74.2 |

| state | StateCount | AverageAge | MaxAge | MinAge |
|-------|------------|------------|--------|--------|
| VA | 8 | 72.83750 | 85.3 | 65.7 |
| OH | 7 | 62.87143 | 72.5 | 49.9 |
| MA | 4 | 72.63333 | 90.7 | 46.5 |
| NY | 4 | 69.32500 | 79.7 | 60.2 |
| TX | 2 | 71.45000 | 78.5 | 64.4 |
| NC | 2 | 60.15000 | 66.6 | 53.7 |

| state | last_name | age |
|-------|-----------|-----|
| MA | Adams | 90.7 |
| VA | Madison | 85.3 |
| SC | Jackson | 78.3 |
| NY | Van Buren | 79.7 |
| NH | Pierce | 64.9 |
| PA | Buchanan | 77.2 |

# Can we do without the views?

```
SELECT
    state,
    last_name as oldest,
    ROUND(DATEDIFF(death, p.birth) / 365, 1) as age

FROM
    sampdb.president p

WHERE
  (   state,
      ROUND(DATEDIFF(death, birth) / 365, 1)
  ) IN
  (
      SELECT
          state,
          MAX(ROUND(DATEDIFF(death, birth) / 365, 1)) as age
      FROM
          sampdb.president p2
      GROUP BY
          state
  )
;
```

- Here we use a "subquery"

- We're still combining tables with one another, but here we generate one of them on the fly