# Practical SQL

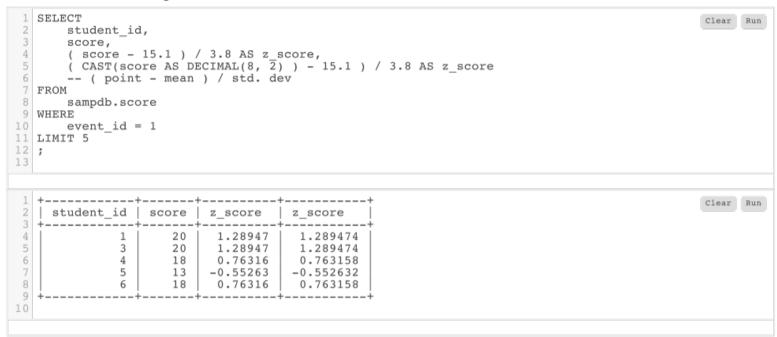# Operations

## We can do simple arithmetic math

```
1  SELECT
2      score / 100 as pct_score,
3      CAST(score AS decimal(8, 2)) / 100 as pct_score_2
4  FROM
5      sampdb.score
6  LIMIT 5;
7
```

Clear  Run

```
1  +-----------+-------------+
2  | pct_score | pct_score_2 |
3  +-----------+-------------+
4  |    0.2000 |    0.200000 |
5  |    0.2000 |    0.200000 |
6  |    0.1800 |    0.180000 |
7  |    0.1300 |    0.130000 |
8  |    0.1800 |    0.180000 |
9  +-----------+-------------+
10
```

Clear  Run

-- Note that MySQL will automatically convert the value returned as a decimal.

-- Other SQL versions will not. Postgres requires a type conversion of the INT to a DECIMAL type.

# More complex math

## We can do more complex math

```
1  SELECT
2      student_id,
3      score,
4      ( score - 15.1 ) / 3.8 AS z_score,
5      ( CAST(score AS DECIMAL(8, 2) ) - 15.1 ) / 3.8 AS z_score
6      -- ( point - mean ) / std. dev
7  FROM
8      sampdb.score
9  WHERE
10     event_id = 1
11 LIMIT 5
12 ;
13
```

Clear   Run

```
1  +------------+-------+----------+-----------+
2  | student_id | score | z_score  | z_score   |
3  +------------+-------+----------+-----------+
4  |          1 |    20 |  1.28947 |  1.289474 |
5  |          3 |    20 |  1.28947 |  1.289474 |
6  |          4 |    18 |  0.76316 |  0.763158 |
7  |          5 |    13 | -0.55263 | -0.552632 |
8  |          6 |    18 |  0.76316 |  0.763158 |
9  +------------+-------+----------+-----------+
10
```

Clear   Run

## We are 'hard coding' the mean and standard deviation values

### We'll discuss how to get these on the fly

# Strings

## String functions manipulate the appearance of text

```
1  SELECT
2      first_name,
3      UPPER(first_name),
4      LOWER(first_name)
5
6  FROM
7      sampdb.president
8
9  LIMIT 5;
10
```

Clear  Run

```
1  +------------+--------------------+--------------------+
2  | first_name | UPPER(first_name)  | LOWER(first_name)  |
3  +------------+--------------------+--------------------+
4  | George     | GEORGE             | george             |
5  | John       | JOHN               | john               |
6  | Thomas     | THOMAS             | thomas             |
7  | James      | JAMES              | james              |
8  | James      | JAMES              | james              |
9  +------------+--------------------+--------------------+
10
```

Clear  Run

# String extraction

## String functions can "reach into" strings

```
 1  SELECT
 2      last_name,
 3      LENGTH(last_name) AS len,
 4      LEFT(last_name, 3) AS lft,
 5      RIGHT(last_name, 3) AS rght,
 6      SUBSTRING(last_name, 3, 2) AS sub,
 7      SUBSTRING(last_name, LENGTH(last_name) - 2, 2) AS penultimate_2
 8  FROM
 9      sampdb.president
10  LIMIT 5;
11
```

Clear  Run

```
 1  +------------+-----+-----+------+-----+---------------+
 2  | last_name  | len | lft | rght | sub | penultimate_2 |
 3  +------------+-----+-----+------+-----+---------------+
 4  | Washington |  10 | Was | ton  | sh  | to            |
 5  | Adams      |   5 | Ada | ams  | am  | am            |
 6  | Jefferson  |   9 | Jef | son  | ff  | so            |
 7  | Madison    |   7 | Mad | son  | di  | so            |
 8  | Monroe     |   6 | Mon | roe  | nr  | ro            |
 9  +------------+-----+-----+------+-----+---------------+
10
```

Clear  Run

# Finding strings allows more interesting usages

## String manipulation becomes more interesting when we start looking for strings

```
1 SELECT
2     email,
3     -- LOCATE('@', email),
4     position('@' in email)
5 FROM
6     sampdb.member
7 LIMIT 5;
8
```
`Clear` `Run`

```
1  +-------------------------+-------------------------+
2  | email                   | position('@' in email)  |
3  +-------------------------+-------------------------+
4  | jeanne_s@earth.com      |                       9 |
5  | august.lundsten@pluto.com |                    16 |
6  | NULL                    |                    NULL |
7  | arbogast.ruth@mars.net  |                      14 |
8  | c.dorfman@uranus.net    |                      10 |
9  +-------------------------+-------------------------+
10
```
`Clear` `Run`

Here we get the position of the @ sign

We can use that, with our other functions, to get a list of domains

# Putting string functions together

## With the location of the '@', we can extract the domain name

```
SELECT
    email,
    RIGHT(email, LENGTH(email) - LOCATE('@', email)) AS domain
FROM
    sampdb.member
LIMIT 5;
```

Clear   Run

```
+--------------------------+------------+
| email                    | domain     |
+--------------------------+------------+
| jeanne_s@earth.com       | earth.com  |
| august.lundsten@pluto.com| pluto.com  |
| NULL                     | NULL       |
| arbogast.ruth@mars.net   | mars.net   |
| c.dorfman@uranus.net     | uranus.net |
+--------------------------+------------+
```

Clear   Run

Now, for each record, we can extract the domain name

Might use this to count members by domain our list use, or check for correlation with other factors

# Date Functions

```
 1  SELECT
 2      birth,
 3      death,
 4      YEAR(birth) AS yr,
 5      MONTH(birth) AS mnt,
 6      MONTHNAME(birth) AS mnt_name,
 7      DATEDIFF(death, birth) AS days,
 8      DATEDIFF(death, birth) / 365 AS years
 9  FROM
10      sampdb.president
11  LIMIT 5;
12
```

```
 1  +------------+------------+------+------+----------+-------+---------+
 2  | birth      | death      | yr   | mnt  | mnt_name | days  | years   |
 3  +------------+------------+------+------+----------+-------+---------+
 4  | 1732-02-22 | 1799-12-14 | 1732 |    2 | February | 24767 | 67.8548 |
 5  | 1735-10-30 | 1826-07-04 | 1735 |   10 | October  | 33119 | 90.7370 |
 6  | 1743-04-13 | 1826-07-04 | 1743 |    4 | April    | 30397 | 83.2795 |
 7  | 1751-03-16 | 1836-06-28 | 1751 |    3 | March    | 31150 | 85.3425 |
 8  | 1758-04-28 | 1831-07-04 | 1758 |    4 | April    | 26729 | 73.2301 |
 9  +------------+------------+------+------+----------+-------+---------+
10
```

# Putting it together

**We can combine records, operations and functions to calculate and format output**

```
 1  SELECT
 2      CONCAT(
 3          first_name,
 4          " ",
 5          last_name,
 6          " lived ",
 7          ROUND(DATEDIFF(death, birth) / 365, 1),
 8          " years"
 9      ) AS sentence
10  FROM
11      sampdb.president
12  LIMIT 5;
13
```

Clear  Run

```
 1  +-----------------------------------+
 2  | sentence                          |
 3  +-----------------------------------+
 4  | George Washington lived 67.9 years |
 5  | John Adams lived 90.7 years        |
 6  | Thomas Jefferson lived 83.3 years  |
 7  | James Madison lived 85.3 years     |
 8  | James Monroe lived 73.2 years      |
 9  +-----------------------------------+
10
```

Clear  Run

# More on Aggregate Functions

# We attack this through dataset exploration

**SELECT DISTINCT gives us the set of values in an attribute -- no repetitions**

```
1  SELECT DISTINCT
2      state
3  FROM
4      sampdb.president
5  ORDER BY
6      state
7  ;
8
```
`Clear`  `Run`

```
1  +-------+
2  | state |
3  +-------+
4  | AR    |
5  | CA    |
6  | CT    |
7  | GA    |
8  | IA    |
9  | IL    |
10 | KY    |
11 | MA    |
12 | MO    |
13 | NC    |
14 | NE    |
15 | NH    |
16 | NJ    |
17 | NY    |
18 | OH    |
19 | PA    |
20 | SC    |
21 | TX    |
22 | VA    |
23 | VT    |
24 +-------+
25
```
`Clear`  `Run`

# This is basically what we get from GROUP BY

## GROUP BY subdivides the dataset into subsets

Each subset is characterized by having the same value for the GROUP BY attribute

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

Clear  Run

```
 1  +-------+-------+
 2  | state | count |
 3  +-------+-------+
 4  | AR    |     1 |
 5  | CA    |     1 |
 6  | CT    |     1 |
 7  | GA    |     1 |
 8  | IA    |     1 |
 9  | IL    |     1 |
10  | KY    |     1 |
11  | MA    |     4 |
12  | MO    |     1 |
13  | NC    |     2 |
14  | NE    |     1 |
15  | NH    |     1 |
16  | NJ    |     1 |
17  | NY    |     4 |
18  | OH    |     7 |
19  | PA    |     1 |
20  | SC    |     1 |
21  | TX    |     2 |
22  | VA    |     8 |
23  | VT    |     2 |
24  +-------+-------+
25
```

Clear  Run

# Same analysis, against member

```
 1  SELECT
 2      state,
 3      COUNT(state) AS count
 4  FROM
 5      sampdb.member
 6  GROUP BY
 7      state
 8  ORDER BY
 9      state
10  LIMIT 15;
11
```

Clear  Run

```
 1  +-------+-------+
 2  | state | count |
 3  +-------+-------+
 4  | AK    |     1 |
 5  | AL    |     3 |
 6  | AZ    |     1 |
 7  | CA    |     6 |
 8  | CO    |     3 |
 9  | CT    |     1 |
10  | FL    |     5 |
11  | GA    |     3 |
12  | HI    |     1 |
13  | IA    |     2 |
14  | ID    |     1 |
15  | IL    |     6 |
16  | IN    |     3 |
17  | KS    |     2 |
18  | KY    |     2 |
19  +-------+-------+
20
```

Clear  Run

# DISTINCT also works by _tuples_

```
 1  SELECT DISTINCT
 2      city,
 3      state
 4  FROM
 5      sampdb.member
 6  ORDER BY
 7      state,
 8      city
 9  LIMIT 15
10  ;
11
```

Clear   Run

```
 1  +---------------+-------+
 2  | city          | state |
 3  +---------------+-------+
 4  | Fairbanks     | AK    |
 5  | Dothan        | AL    |
 6  | Huntsville    | AL    |
 7  | Mobile        | AL    |
 8  | Kayenta       | AZ    |
 9  | Los Angeles   | CA    |
10  | Oakland       | CA    |
11  | San Francisco | CA    |
12  | Stockton      | CA    |
13  | Trona         | CA    |
14  | Denver        | CO    |
15  | Durango       | CO    |
16  | Waterbury     | CT    |
17  | Coral Gables  | FL    |
18  | Fort Myers    | FL    |
19  +---------------+-------+
20
```

Clear   Run

Here we are getting a list of city / state value pairs that occur at least once in the data

This is not _permutations_, e.g. the possible combinations of these values -- each of these _does_ appear at least once

# As does GROUP BY

## Grouping by _several_ attributes subdivides counts

```
 1  SELECT
 2      city,
 3      state,
 4      COUNT(*)
 5  FROM
 6      sampdb.member
 7  GROUP BY
 8      state,
 9      city
10  ORDER BY
11      state,
12      city
13  LIMIT 15
14  ;
15
```

`Clear` `Run`

```
 1  +---------------+-------+----------+
 2  | city          | state | COUNT(*) |
 3  +---------------+-------+----------+
 4  | Fairbanks     | AK    |        1 |
 5  | Dothan        | AL    |        1 |
 6  | Huntsville    | AL    |        1 |
 7  | Mobile        | AL    |        1 |
 8  | Kayenta       | AZ    |        1 |
 9  | Los Angeles   | CA    |        2 |
10  | Oakland       | CA    |        1 |
11  | San Francisco | CA    |        1 |
12  | Stockton      | CA    |        1 |
13  | Trona         | CA    |        1 |
14  | Denver        | CO    |        2 |
15  | Durango       | CO    |        1 |
16  | Waterbury     | CT    |        1 |
17  | Coral Gables  | FL    |        1 |
18  | Fort Myers    | FL    |        1 |
19  +---------------+-------+----------+
20
```

`Clear` `Run`

All the Alabama records are still in the output.

But the sum of Alabama records, which was 3, is now subdivided over three different cities

# ORDER BY shows the grouping of the records

```
 1 SELECT
 2     last_name,
 3     first_name,
 4     member_id,
 5     city,
 6     state
 7 FROM
 8     sampdb.member
 9 ORDER BY
10     state,
11     city
12 LIMIT 15
13 ;
14
```

Clear  Run

```
 1 +-----------+------------+-----------+----------------+-------+
 2 | last_name | first_name | member_id | city           | state |
 3 +-----------+------------+-----------+----------------+-------+
 4 | Matthews  | Bill       |        56 | Fairbanks      | AK    |
 5 | Edwards   | John       |        82 | Dothan         | AL    |
 6 | Hughes    | Max        |        37 | Huntsville     | AL    |
 7 | Schauer   | Alma       |        65 | Mobile         | AL    |
 8 | Kirby     | Timothy    |        48 | Kayenta        | AZ    |
 9 | Puntillo  | Cheryl     |        59 | Los Angeles    | CA    |
10 | Pierson   | Stanley    |        88 | Los Angeles    | CA    |
11 | Sprague   | Earl       |        99 | Oakland        | CA    |
12 | Smith     | Laura      |        98 | San Francisco  | CA    |
13 | Feit      | Daniel     |        19 | Stockton       | CA    |
14 | Simmons   | David      |        49 | Trona          | CA    |
15 | Garner    | Steve      |        89 | Denver         | CO    |
16 | Sawyer    | Dennis     |         7 | Denver         | CO    |
17 | Bookstaff | Barbara    |        47 | Durango        | CO    |
18 | Schenk    | Cindy      |        42 | Waterbury      | CT    |
19 +-----------+------------+-----------+----------------+-------+
20
```

Clear  Run

# Referring to SQL results in SQL queries -- Subqueries

# Subqueries allow us to 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

# Presidents from the most presidential states

## 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
```

Clear   Run

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

Clear   Run

# Hard-code the values

## We can put those values into a query

```
 1  SELECT
 2      last_name,
 3      first_name,
 4      state
 5  FROM
 6      sampdb.president
 7  WHERE
 8      state IN
 9      (
10          'VA',
11          'OH',
12          'MA'
13      )
14  ORDER BY
15      state
16  ;
17
18
```

<span style="float:right">Clear Run</span>

```
 1  +------------+----------------+-------+
 2  | last_name  | first_name     | state |
 3  +------------+----------------+-------+
 4  | Bush       | George H.W.    | MA    |
 5  | Adams      | John Quincy    | MA    |
 6  | Kennedy    | John F.        | MA    |
 7  | Adams      | John           | MA    |
 8  | Harding    | Warren G.      | OH    |
 9  | Taft       | William H.     | OH    |
10  | McKinley   | William        | OH    |
11  | Harrison   | Benjamin       | OH    |
12  | Garfield   | James A.       | OH    |
13  | Hayes      | Rutherford B.  | OH    |
14  | Grant      | Ulysses S.     | OH    |
15  | Taylor     | Zachary        | VA    |
16  | Tyler      | John           | VA    |
17  | Harrison   | William H.     | VA    |
18  | Monroe     | James          | VA    |
19  | Wilson     | Woodrow        | VA    |
20  | Madison    | James          | VA    |
21  | Jefferson  | Thomas         | VA    |
22  | Washington | George         | VA    |
23  +------------+----------------+-------+
24
```

<span style="float:right">Clear Run</span>

## But, Kludgy!

What if the data changes?

# Joins

# A basic query:

SELECT
   student_id AS Student,
   AVG(score) AS Average,
   COUNT(score) AS "# of Tests"
FROM
   sampdb.score
GROUP BY
   student_id
ORDER BY
   Average DESC

- Returns each student's average

| Student | Average | # of Tests |
|---------|---------|-----------|
| ▶ 1 | 48.0000 | 5 |
| 27 | 45.7500 | 4 |
| 5 | 42.8333 | 6 |
| 18 | 41.3333 | 6 |
| 17 | 41.2000 | 5 |
| 2 | 40.4000 | 5 |
| 11 | 39.8333 | 6 |

- But can we get the student's name rather than their id?

# Yes

```
SELECT
    st.name AS        Name,
    scr.student_id AS   Id,
    AVG(scr.score) AS   Average,
    COUNT(scr.score) AS "# Tests"
FROM
    sampdb.score scr
INNER JOIN
    sampdb.student st
ON
    scr.student_id = st.student_id
GROUP BY
    scr.student_id
ORDER BY
    Average DESC
```

| Name | Id | Average | # Tests |
|------|----|---------|---------|
| Megan | 1 | 48.0000 | 5 |
| Carter | 27 | 45.7500 | 4 |
| Abby | 5 | 42.8333 | 6 |
| Max | 18 | 41.3333 | 6 |
| Will | 17 | 41.2000 | 5 |
| Joseph | 2 | 40.4000 | 5 |

# Breaking down the query

```
SELECT
    st.name AS        Name,
    scr.student_id AS   Id,
    AVG(scr.score) AS   Average,
    COUNT(scr.score) AS "# Tests"
FROM
    sampdb.score scr
INNER JOIN
    sampdb.student st
ON
    scr.student_id = st.student_id
GROUP BY
    scr.student_id
ORDER BY
    Average DESC
```

- A join combines rows from different tables
- The "ON" clause specifies conditions for combining records
  - Here, records with the same student_id value are combined into a single record
  - The SELECT clause still specifies which fields are displayed from that combined record
- INNER JOIN specifies which table to join data from
  - The addition of "scr" and "st" in the specification of tables provides alias values for reference by the other clauses
  - Without these the query wouldn't know which table it should find a field in

# Two Really Simple Tables

- T1

```
+-------+-------+
| i1    | c1    |
+-------+-------+
|     1 | a     |
|     2 | b     |
|     3 | c     |
+-------+-------+
```

- t2

```
+-------+-------+
| i2    | c2    |
+-------+-------+
|     2 | c     |
|     3 | b     |
|     4 | a     |
+-------+-------+
```

# Inner Join: All Rows Matched

- A simple inner join matches each row in one table with each row in the other

- So joining a 3 row table with a 3 row table produces a 9 row table

- Obviously we don't want all those rows

SELECT * FROM join_sample.t1 INNER JOIN join_sample.t2;

```
+------+------+------+------+
| i1   | c1   | i2   | c2   |
+------+------+------+------+
|    1 | a    |    2 | c    |
|    2 | b    |    2 | c    |
|    3 | c    |    2 | c    |
|    1 | a    |    3 | b    |
|    2 | b    |    3 | b    |
|    3 | c    |    3 | b    |
|    1 | a    |    4 | a    |
|    2 | b    |    4 | a    |
|    3 | c    |    4 | a    |
+------+------+------+------+
```

# Limiting Inner Join Results

SELECT
  *

FROM
  join_sample.t1

INNER JOIN
  join_sample.t2

WHERE
  join_sample.t1.i1 = join_sample.t2.i2

;

```
+-------+-------+-------+-------+
|  i1   |  c1   |  i2   |  c2   |
+-------+-------+-------+-------+
|   1   | a     |   2   | c     |
|   2   | b     |   2   | c     |
|   3   | c     |   2   | c     |
|   1   | a     |   3   | b     |
|   2   | b     |   3   | b     |
|   3   | c     |   3   | b     |
|   1   | a     |   4   | a     |
|   2   | b     |   4   | a     |
|   3   | c     |   4   | a     |
+-------+-------+-------+-------+
```

```
+-------+-------+-------+-------+
|  i1   |  c1   |  i2   |  c2   |
+-------+-------+-------+-------+
|   2   | b     |   2   | c     |
|   3   | c     |   3   | b     |
+-------+-------+-------+-------+
```

# Left Join

- Get all rows from the "left" table, each with that row from the "right" table that matches on the specified fields

```
SELECT
  *
FROM
  join_sample.t1
LEFT JOIN
  join_sample.t2
  ON
    join_sample.t1.i1 = join_sample.t2.i2
;
```

```
+-------+-------+-------+-------+
| i1    | c1    | i2    | c2    |
+-------+-------+-------+-------+
|     1 | a     |  NULL | NULL  |
|     2 | b     |     2 | c     |
|     3 | c     |     3 | b     |
+-------+-------+-------+-------+
```

# Contact

gordon@practicalhorseshoeing.com