Introduction

POSTGRESQL SUMMARY STATS AND WINDOW FUNCTIONS



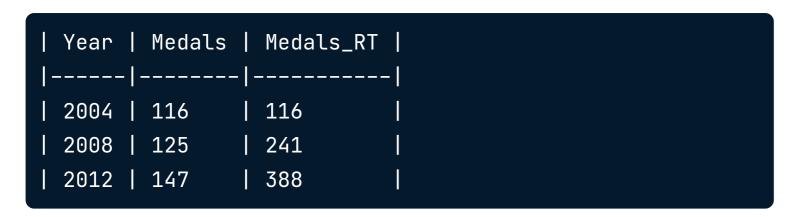
Michel Semaan

Data Scientist



Motivation

USA total and running total of Summer Olympics gold medals since 2004



Discus throw reigning champion status

```
Year | Champion | Last_Champion | Reigning_Champion |
1996 | GER
               | null
                               | false
               l GER
                               | false
               | LTU
                               | true
2008
    l EST
               | LTU
                               | false
2012
      GER
               | EST
                               | false
```

Course outline

- 1. Introduction to window functions
- 2. Fetching, ranking, and paging
- 3. Aggregate window functions and frames
- 4. Beyond window functions

Summer olympics dataset

Each row represents a medal awarded in the Summer Olympics games

Columns

- Year, City
- Sport, Discipline, Event
- Athlete, Country, Gender
- Medal

Window functions

- Perform an operation across a set of rows that are somehow related to the current row
- Similar to GROUP BY aggregate functions, but all rows remain in the output

Uses

- Fetching values from preceding or following rows (e.g. fetching the previous row's value)
 - Determining reigning champion status
 - Calculating growth over time
- Assigning ordinal ranks (1rst, 2nd, etc.) to rows based on their values' positions in a sorted list
- Running totals, moving averages



Row numbers

Query

```
SELECT
  Year, Event, Country
FROM Summer_Medals
WHERE
  Medal = 'Gold';
```

Result

Enter ROW_NUMBER

Query

```
SELECT
  Year, Event, Country,
  ROW_NUMBER() OVER () AS Row_N
FROM Summer_Medals
WHERE
  Medal = 'Gold';
```

Result

Anatomy of a window function

Query

```
SELECT
  Year, Event, Country,
  ROW_NUMBER() OVER () AS Row_N
FROM Summer_Medals
WHERE
  Medal = 'Gold';
```

- FUNCTION_NAME() OVER (...)
 - ORDER BY
 - PARTITION BY
 - ROWS/RANGE PRECEDING/FOLLOWING/UNBOUNDED

Let's practice!

POSTGRESQL SUMMARY STATS AND WINDOW FUNCTIONS



ORDER BY

POSTGRESQL SUMMARY STATS AND WINDOW FUNCTIONS



Michel Semaan

Data Scientist



Row numbers

Query

```
SELECT
  Year, Event, Country,
  ROW_NUMBER() OVER () AS Row_N
FROM Summer_Medals
WHERE
  Medal = 'Gold';
```

Result*

Enter ORDER BY

- ORDER BY in OVER orders the rows related to the current row
 - Example: Ordering by year in descending order in ROW_NUMBER 's OVER clause will assign
 1 to the most recent year's rows

Ordering by Year in descending order

Query Result

```
SELECT
  Year, Event, Country,
  ROW_NUMBER() OVER (ORDER BY Year DESC) AS Row_N
FROM Summer_Medals
WHERE
  Medal = 'Gold';
```

Ordering by multiple columns

Query Result

```
SELECT
  Year, Event, Country,
  ROW_NUMBER() OVER
    (ORDER BY Year DESC, Event ASC) AS Row_N
FROM Summer_Medals
WHERE
  Medal = 'Gold';
```

```
Year | Event
               | Country | Row_N |
2012 | + 100KG
              | FRA
2012 | + 67 KG |
                 SRB
2012 | + 78KG |
                 CUB
```

Ordering in- and outside OVER

Query

```
SELECT
  Year, Event, Country,
  ROW_NUMBER() OVER
    (ORDER BY Year DESC, Event ASC) AS Row_N
FROM Summer_Medals
WHERE
  Medal = 'Gold'
ORDER BY Country ASC, Row_N ASC;
```

Result

ORDER BY inside OVER takes effect before
 ORDER BY outside OVER

Reigning champion

- A reigning champion is a champion who's won both the previous and current years' competitions
- The previous and current year's champions need to be in the same row (in two different columns)

Enter LAG

- LAG(column, n) OVER (...) returns column 's value at the row n rows before the current row
 - LAG(column, 1) OVER (...) returns the previous row's value

Current champions

Query

```
SELECT
  Year, Country AS Champion
FROM Summer_Medals
WHERE
  Year IN (1996, 2000, 2004, 2008, 2012)
  AND Gender = 'Men' AND Medal = 'Gold'
  AND Event = 'Discus Throw';
```

Result

Current and last champions

Query

```
WITH Discus_Gold AS (
  SELECT
    Year, Country AS Champion
  FROM Summer Medals
  WHERE
    Year IN (1996, 2000, 2004, 2008, 2012)
    AND Gender = 'Men' AND Medal = 'Gold'
    AND Event = 'Discus Throw')
SELECT
  Year, Champion,
  LAG(Champion, 1) OVER
    (ORDER BY Year ASC) AS Last_Champion
FROM Discus_Gold
ORDER BY Year ASC;
```

Result

Let's practice!

POSTGRESQL SUMMARY STATS AND WINDOW FUNCTIONS



PARTITION BY

POSTGRESQL SUMMARY STATS AND WINDOW FUNCTIONS



Michel Semaan

Data Scientist



Motivation

Query

```
WITH Discus_Gold AS (
  SELECT
    Year, Event, Country AS Champion
  FROM Summer Medals
  WHERE
    Year IN (2004, 2008, 2012)
    AND Gender = 'Men' AND Medal = 'Gold'
    AND Event IN ('Discus Throw', 'Triple Jump')
    AND Gender = 'Men')
SELECT
  Year, Event, Champion,
  LAG(Champion) OVER
    (ORDER BY Event ASC, Year ASC) AS Last_Champion
FROM Discus_Gold
ORDER BY Event ASC, Year ASC;
```

Result

```
| Champion | Last_Champion |
Year | Event
2004 | Discus Throw |
                     LTU
                               I null
      Discus Throw | EST
                               l LTU
      Discus Throw | GER
                                EST
                                GER
2004 | Triple Jump
                   l SWE
2008 | Triple Jump
                   I POR
                               I SWE
2012 | Triple Jump | USA
                                POR
```

 When Event changes from Discus Throw to Triple Jump, LAG fetched
 Discus Throw 's last champion as opposed to a null

Enter PARTITION BY

- PARTITION BY splits the table into partitions based on a column's unique values
 - The results aren't rolled into one column
- Operated on separately by the window function
 - ROW_NUMBER will reset for each partition
 - LAG will only fetch a row's previous value if its previous row is in the same partition

Partitioning by one column

Query

```
WITH Discus_Gold AS (...)

SELECT
   Year, Event, Champion,
   LAG(Champion) OVER
        (PARTITION BY Event
        ORDER BY Event ASC, Year ASC) AS Last_Champion
FROM Discus_Gold
ORDER BY Event ASC, Year ASC;
```

Result

```
| Champion | Last_Champion |
Year | Event
2004 | Discus Throw |
                     LTU
                               | null
      Discus Throw |
                     EST
                               l LTU
      Discus Throw |
                                 EST
                     GER
2004 | Triple Jump
                               I null
                   I SWE
    | Triple Jump
                     POR
                                SWE
2012 | Triple Jump |
                     USA
                                 POR
```



More complex partitioning

```
Year | Country | Event
                                   | Row_N |
              | + 78KG (Heavyweight) | 1
2008 | CHN
              | - 49 KG
              | 48 - 55KG
2008 | JPN
           | 48 - 55KG
2008 | JPN
              +75KG
                                   | 32
2012 | CHN
2012 | CHN
            | - 49 KG
              +75KG
2012 | JPN
           | - 49 KG
2012 | JPN
              . . . .
```

Row number should reset per Year and Country

Partitioning by multiple columns

Query Result

```
WITH Country_Gold AS (
  SELECT
    DISTINCT Year, Country, Event
  FROM Summer Medals
  WHERE
    Year IN (2008, 2012)
    AND Country IN ('CHN', 'JPN')
    AND Gender = 'Women' AND Medal = 'Gold')
SELECT
  Year, Country, Event,
  ROW_NUMBER() OVER (PARTITION BY Year, Country)
FROM Country_Gold;
```

```
Year | Country | Event
                                     Row_N
              | + 78KG (Heavyweight) | 1
2008 | CHN
              l - 49 KG
2008
      CHN
              | ...
              | 48 - 55KG
2008
    l JPN
2008
     l JPN
              | 48 - 55KG
2012
      CHN
              l +75KG
2012
      CHN
                - 49 KG
               +75KG
2012
    JPN
2012 | JPN
               - 49 KG
```

Let's practice!

POSTGRESQL SUMMARY STATS AND WINDOW FUNCTIONS

