

# Appunti Basi di Dati

Luca Seggiani

23 Maggio 2024

## 1 Data analytics

### Window function

Una window function (ovvero funzione analitica) è una funzione che associa ad ogni record  $r$  un valore ottenuto da un'operazione su insieme di record logicamente connessi ad  $r$ , come ad esempio media, rank, ecc... Vediamo un'esempio di applicazione delle window function. Vogliamo scrivere una query che indichi, per ogni cardiologo, la matricola, la parcella, e la parcella media della sua specializzazione. Senza usare correlated subquery nel select, si può invece usare la clausola OVER. La clausola OVER applica una funzione di tipo aggregato o non-aggregato a un'insieme di record associati a un record di result set, detto *partition*. L'esempio sarebbe quindi:

```
1 SELECT M.Matricola ,  
2   M.Parcella ,  
3   M.Parcella  
4     AVG(M.Parcella) OVER()  
5 FROM Medico M  
6 WHERE M.Specializzazione = 'Cardiologia';
```

Si può generalizzare questa query a tutte le specializzazioni attraverso la partizione. La partizione è l'insieme di record a cui si applica la funzione aggregata / non aggregata, e dipende dal record processato attualmente. L'esempio può essere:

```
1 SELECT M.Matricola ,  
2   M.Parcella ,  
3   M.Specializzazione ,  
4   M.Parcella  
5     AVG(M.Parcella) OVER( PARTITION BY  
6                               M.Specializzazione)  
7 FROM Medico M
```

Dove la clausola PARTITION BY M.Specializzazione richiede alla query di effettuare l'operazione di aggregazione AVG su sottoinsiemi formati da una partizione della tabella Medico contenente solo medici con la stessa specializzazione del medico d'interesse. Il funzionamento è simile a quello di una GROUP BY, o di una semplice aggregazione, ma mantiene l'informazione riguardante i singoli medici (da cui ricaviamo matricola e parcella). Con OVER si possono usare diverse funzioni aggregate, una cui lista non viene riportata qua. Di queste abbiamo già visto le AVG, COUNT, COUNT(DISTINCT), MAX, MIN e SUM. Notiamo poi che la definizione di WINDOW è concessa fuori dalla OVER:

```

1 SELECT M.Matricola ,
2     AVG(...) OVER w
3 FROM Medico M
4 WINDOW w AS ( ... )

```

### Funzioni non aggregate

Le funzioni non aggregate si possono dividere in due sottocategorie: quelle che usano l'intera partizione, e quelle che usano un suo sottoinsieme (*frame*). Le window functions non aggregate associano a ciascun record di un result set un valore ottenuto senza l'aggregazione del result set stesso, ma calcolato su tutti i record a partire dall'inizio della partizione fino al record attuale se si specifica un ORDER BY nella partizione, e tutti i record della partizione in caso contrario. Alcuni esempi sono ad esempio la FIRST\_VALUE, LAST\_VALUE, ecc... Possiamo fare un'esempio: poniamo di voler associare un numero ad ogni medico nella sua specializzazione. Potremo usare la funzione ROW\_NUMBER, che restituisce un'indice all'interno della partizione per ogni record:

```

1 SELECT M.Matricola ,
2     M.Cognome ,
3     M.Specializzazione ,
4     ROW_NUMBER() OVER( PARTITION BY
5                         M.Specializzazione)
6 FROM Medico M;

```

### Classificazione

Vediamo l'uso della funzione RANK per stilare classifiche. La funzione RANK permette di classificare record (banalmente, ordinare) associando ad ognuno un certo valore SCORE (punteggio). Poniamo per esempio di voler classificare i medici in base alla loro convenienza (quindi in base a parcelle più basse):

```

1 SELECT M.Matricola ,
2     M.Cognome ,
3     M.Specializzazione ,
4     M.Parcella

```

```

5 RANK() OVER( ORDER BY M.Parcella )
6 FROM Medico M;

```

Questo query, per ogni medico, ottiene il sottoinsieme della partizione del resultset ordinato fino a quel medico (con ORDER BY), e restituisce la posizione (con RANK) di quel medico in tale partizione. Di default, la ORDER BY effettua un'ordinamento ascendente (quindi parcelle più basse in cima, parcelle più alte in fondo). Vediamo la gestione dei pareggi. In questa forma, la funzione rank salta, per ogni parimerito, un numero uguale al numero dei pareggi, di posizioni. Ergo, se due medici ottengono lo stesso punteggio, si avrà:

Nome	Parcella	Posizione
Fabio Rossi	150	1
Nicola Tonini	250	2
Tullio Nomefalso	250	2
Michele Poraccio	300	4

Per ottenere un comportamento diverso posso usare il DENSE\_RANK. Il DENSE\_RANK non scala la posizione nel caso di un parimerito, e dà quindi una classifica senza "buchi":

```

1 SELECT M.Matricola ,
2   M.Cognome ,
3   M.Specializzazione ,
4   M.Parcellam
5   DENSE_RANK() OVER( ORDER BY M.Parcella )
6 FROM Medico M;

```

Nome	Parcella	Posizione
Fabio Rossi	150	1
Nicola Tonini	250	2
Tullio Nomefalso	250	2
Michele Poraccio	300	3

Si possono effettuare classifiche anche in base su più sottinsiemi (i cosiddetti **rank multipli**). Stiliamo, ad esempio, la classifica delle visite effettuate sia riguardo alla totalità dei medici, sia riguardo alla specializzazione di appartenenza del medico:

```

1 WITH visite AS (
2   SELECT V.Medico ,
3     M.Cognome ,
4     M.Specializzazione ,
5     COUNT(*) AS Visite

```

```

6   FROM Visita V
7     INNER JOIN Medico M ON M.Matricola = V.Medico
8   GROUP BY V.Medico
9 )
10  SELECT VV.Cognome,
11    VV.Specializzazione,
12    VV.Visite,
13    RANK() OVER(ORDER BY VV.Visite DESC) AS GlobalRank ,
14    RANK() OVER(PARTITION BY VV.Specializzazione
15                      ORDER BY VV.Visite DESC) AS SpecRank
16  FROM visite VV;

```

## Funzioni LEAD e LAG

Le funzioni LEAD e LAG ricavano il valore dell'attributo (che può essere anche ricavato) di una row posta k posizioni prima o dopo un dato record. Chiaramente, per applicare le LEAD e LAG dev'essere stabilito una relazione di ordinamento sensata nella partizione considerata. La funzione LAG, ad esempio, prende due argomenti: il primo rappresenta l'attributo su cui effettuare il "salto", mentre il secondo rappresenta il numero di posizioni da saltare. Poniamo ad esempio di voler considerare, per tutte le visite otorinolaringoiatriche dal 2010 al 2019, la matricola del medico, il codice fiscale del paziente, la data e la data della visita precedente effettuata da un medico con la stessa specializzazione. Avremo:

```

1  SELECT V.Medico ,
2    V.Paziente ,
3    V.Data ,
4    LAG(V.Data , 1) OVER ( PARTITION BY V.Paziente
5                           ORDER BY V.Data)
6  FROM Visita V
7    INNER JOIN Medico M ON V.Medico = M.Matricola
8 WHERE M.Specializzazione = "Otorinolaringoiatria"
9   AND YEAR(V.Data) BEWEEN 2010 AND 2019

```

## Frame

Un frame è un sottoinsieme di una partizione. Per la precisione, a partire dal record corrente in una partizione, ci si può muovere N posizioni **preceding** (precedenti) e N posizioni **following** (seguenti) da considerare. I record al di fuori delle soglie preceding e following si chiamano preceding o following **unbounded**. Una window function su un frame processa un result set e calcola valori aggregati e non sulla base di record adiacenti alla current row, ovvero su un frame della current row. Esistono alcune funzioni aggregate che lavorano solo su frame: ad esempio FIRST\_VALUE e LAST\_VALUE, che restituiscono rispettivamente il primo e l'ultimo valore del frame. Se non si specifica un frame, viene impostato automaticamente un default frame (che non corrisponde necessariamente alla partizione intera). Il default frame

equivale al frame che contiene tutti valori precedenti nel caso di OVER con ORDER BY, e all'intera partizione nel caso di OVER semplici. Vediamo un'esempio dell'uso di FIRST\_VALUE. Poniamo di voler ottenere le prime visite cardiologiche fatte da 3 specifici pazienti nel triennio 2012-2014 con ciascun medico:

```
1 SELECT V.Medico,
2     V.Paziente,
3     V.Data,
4     FIRST_VALUE(V.Data) OVER w
5 FROM Visita V
6 WHERE V.Paziente IN ('aaa_1', 'bbc_4', 'ccc_2')
7     AND YEAR(V.Data) BETWEEN 2012 AND 2014
8 WINDOW w AS (PARTITION BY V.Medico, V.Paziente
9                 ORDER BY V.Data)
```

In questo caso il default frame è corretto, in quanto vogliamo il primo record del result set.