

Window value | Analytics Functions

These functions return values from other rows in the same window.

lag()

lead()

first_value()

last_value()

nth_value()

lag():

used to access a value from a previous row in the same window.

Mandatory requirements:

↳ use lag(): → need must \Rightarrow 1) over() 2) orderBy()

{without order; 'previous row' has no meaning}

① Basic lag() &:

```
df.withColumn(
    'prev_salary',
    lag('salary').over(
        Window.orderBy('salary')
    )
)
```

| Emp | Salary | prev_salary |
|-----|--------|-------------|
| D | 900 | NULL |
| A | 1000 | 900 |
| B | 1100 | 1000 |
| C | 1500 | 1100 |
| | 1500 | 1500 |

(first row has no previous ~~row~~ \Rightarrow NULL)

② lag() with partitionBy:

```
df.withColumn(
    'prev_salary',
    lag('salary').over(
        Window.partitionBy('dept')
    )
)
```

| EMP | dept | salary | prev_salary |
|-----|------|--------|-------------|
| A | IT | 1000 | NULL |
| B | IT | 1500 | 1000 |
| C | IT | 1500 | 1500 |
| D | HR | 900 | 1500 |
| E | HR | 1200 | 900 |

③ lag(column, offset) - jump back more rows

lag('salary', 2)

gives me the salary from 2 rows before

```
df.withColumn(  
    'sal-2-back',  
    lag('salary', 2).over(  
        Window.orderBy('salary'))  
)
```

dp

| emp | salary | sal-2-back |
|-----|--------|------------|
| D | 900 | NULL |
| A | 1000 | NULL |
| E | 1100 | 900 |
| B | 1500 | 1000 |
| C | 1500 | 1100 |

④ lag(column, offset, default):

lag('salary', 1, 0)

{ if previous row doesn't exist, return 0 instead
 ↳ null }

```
df.withColumn(  
    'salary-prev',  
    lag('salary', 1, 0).over(  
        Window.partitionBy('dept')  
            .orderBy('salary'))  
)
```

dp

| dept | emp | salary | prev_salary |
|------|-----|--------|-------------|
| IT | A | 1000 | 0 |
| IT | B | 1500 | 1000 |
| IT | C | 1500 | 1500 |
| HR | D | 900 | 0 |
| HR | E | 1100 | 900 |

③ most common use case - (difference calculation)
↓
(salary increase)

df.withColumn()

'salary-diff'

col('salary') - lag('salary').over(

Window.orderBy('salary'))

| salary | salary-diff |
|--------|-------------|
| 900 | NULL |
| 1000 | 100 |
| 1100 | 100 |
| 1500 | 400 |
| 1500 | 0 |

Real-world production use cases:

- Day-over-day sales diff.
- Month-over-month growth
- Previous status comparison
- Trend detection

Note:

→ `lag()`, `lead()` ignore frame clauses
→ They only depend on `order` and `partition`

How does lag() / lead() behave when a frame clause is present?

↳ lag() / lead() will completely IGNORES the frame clause

↳ Even Default frame also be ignored

Why?

Because lag() / lead() is NOT an aggregate function

Therefore;

Frame clause has no effect on lag() / lead()

lead()

calculator - 10/10

↳ lead() is exactly opposite to lag()
↳ used to access a value from a next row
in the same window.

→ To use lead(); over() orderBy() are mandatory.

→ frame clause has no effect on lead()

① Basic lead():

df.withColumn(
'next-salary',
lead('salary').over(
window.orderBy('salary')))

| | salary | next-salary |
|--|--------|-------------|
| | 900 | 1000 |
| | 1000 | 1100 |
| | 1100 | 1500 |
| | 1500 | 1500 |
| | 1500 | NULL |

② lead() with partitionBy:

df.withColumn(
'next-sal',
lead('salary', 2, 0).over(
window.partitionBy('dept')
.orderBy('salary')))

| emp | dept | salary | next-salary |
|-----|------|--------|-------------|
| A | IT | 1000 | 1500 |
| B | IT | 1500 | 0 |
| C | IT | 1500 | 0 |
| D | HR | 900 | 0 |
| E | HR | 1100 | 0 |

Real world Production Use cases:

- > Next-day sales comparison
- > Event sequence Analysis
- > Time-to-next-event calculation.

first_value()

This function looks simple, but it's one of the most confusing in practice because of the frame clause. Frame clause is implemented here.

first_value(): return the FirstValue in the window frame for the current row

→ mandatory field)

→ orderBy: decides order

frame clause: decides which rows are visible

firstvalues: picks the first row from the visible set

| Emp_id | Emp_name | dept | month | salary |
|--------|----------|------|-------|--------|
| 101 | A | IT | Jan | 50000 |
| 101 | A | IT | Feb | 52000 |
| 101 | A | IT | Mar | 55000 |
| 102 | B | IT | Jan | 60000 |
| 102 | B | IT | Mar | 65000 |
| 103 | C | HR | Jan | 40000 |
| 103 | C | HR | Feb | 42000 |
| 103 | C | HR | Mar | 45000 |

① first_value() with Default frame:

df.withColumn(

'starting_salary',

first_value('salary').over(

window.partitionBy('emp_id')

• orderBy ('~~salary~~' month)

| emp_id | month | salary | sliding_sales |
|--------|-------|--------|---------------|
| 101 | Jan | 30000 | 30K |
| 101 | Feb | 32K | 50K |
| 101 | Mar | 35K | 50K |
| 102 | Jan | 60K | 60K |
| 102 | Mar | 65K | 60K |
| 103 | Jan | 40K | 40K |
| 103 | Feb | 42K | 40K |
| 103 | Mar | 45K | 40K |

Default frame:

(ROWS_BETWEEN UNBOUNDED_PRECEDING AND CURRENT_ROW)

Examples

① Case 2: Full partition frame.

df.withColumn(

'salary-first',

first_value('salary').over(

Window.partitionBy('emp_id')

.orderBy('month')

• RowsBetween(

Window.unboundedPreceding,

Window.unboundedFollowing

| emp_id | month | salary | salary-first |
|--------|-------|--------|--------------|
| 101 | Jan | 30K | 30K |
| 101 | Feb | 32K | 50K |
| 101 | Mar | 35K | 50K |
| 102 | Jan | 60K | 60K |
| 102 | Mar | 65K | 60K |
| 103 | Jan | 40K | 40K |
| 103 | Feb | 42K | 40K |
| 103 | Mar | 45K | 40K |

Case3: current row only

```

df.withColumn(
    'salary_start',
    first_value('salary').over(
        Window.partitionBy('emp_id')
        .start.orderBy('month')
        .rowsBetween(0,0)
    )
)
  
```

| emp_id | salary | month | salary_start |
|--------|--------|-------|--------------|
| 101 | 30k | Jan | 30k |
| 101 | 32k | Feb | 32k |
| 101 | 35k | Mar | 35k |
| 102 | 60k | Jan | 60k |
| 102 | 65k | Mar | 65k |
| 103 | 40k | Jan | 40k |
| 103 | 42k | Feb | 42k |
| 103 | 45k | Mar | 45k |

Case4: current row → Future row

```

df.withColumn(
    'start_salary',
    first_value('salary').over(
        Window.partitionBy('emp_id')
        .orderBy('month')
        .rowsBetween(0,
            Window.unboundedFollowing
        )
    )
)
  
```

| emp-id | salary | month | salary_start |
|--------|--------|-------|--------------|
| 101 | 30k | Jan | 30k |
| 101 | 52k | Feb | 52k |
| 101 | 55k | Mar | 55k |
| 102 | 60k | Jan | 60k |
| 102 | 65k | Mar | 65k |
| 103 | 40k | Jan | 40k |
| 103 | 42k | Feb | 42k |
| 103 | 45k | Mar | 45k |

case5: Sliding window (previous + current)

df.withColumn(

'salary_start',

first_value('salary').over(

window.partitionBy('emp-id').

• orderBy('month')

• rowBetween(-1, 0)

)

| emp-id | salary | month | salary_start |
|--------|--------|-------|--------------|
| 101 | 30k | Jan | 30k |
| 101 | 52k | Feb | 52k |
| 101 | 55k | Mar | 55k |
| 102 | 60k | Jan | 60k |
| 102 | 65k | Mar | 65k |
| 103 | 40k | Jan | 40k |
| 103 | 42k | Feb | 42k |
| 103 | 45k | Mar | 45k |

Gold rule:

first_value() returns first row in the frame

Real world production use cases:

- Data normalization

- compare each row against the starting value

- Baseline comparison

- first sale of the day

- opening price

- initial state

last_value():

Returns the last-value in the window PRAMA for the current row.

Case1: Default ~~row~~ frame:

\downarrow
rowBetween(unboundedPreceding, current row)

`last_value('salary').over(
 Window.partitionBy('emp-id').orderBy('month')`

| emp-id | month | salary | last_sal |
|--------|-------|--------|----------|
| 101 | Jan | 50k | 50k |
| 101 | Feb | 52k | 52k |
| 101 | Mar | 55k | 55k |
| 102 | Jan | 60k | 60k |
| 102 | Mar | 65k | 65k |
| 103 | Jun | 40k | 40k |
| 103 | Feb | 42k | 42k |
| 103 | May | 45k | 45k |

Case2: full partition frame:

`last_value('salary').over(
 Window.partitionBy('emp-id')`

\cdot orderBy('month')
 \cdot rowBetween (

Window.unboundedPreceding,
Window.unboundedFollowing

| emp-id | month | salary | last_sal |
|--------|-------|--------|----------|
| 101 | Jan | 50k | 55k |
| 101 | Feb | 52k | 55k |
| 101 | Mar | 55k | 55k |
| 102 | Jan | 60k | 65k |
| 102 | Mar | 65k | 65k |
| 103 | Jan | 40k | 45k |
| 103 | Feb | 42k | 45k |
| 103 | May | 45k | 45k |

Case3: current row:

`last_value('salary').over(
 Window.partitionBy('emp-id')
 .orderBy('month')
 .rowBetween(0,0))`

| emp-id | month | salary | last_sal |
|--------|-------|--------|----------|
| 101 | Jan | 50k | 50k |
| 101 | Feb | 52k | 52k |
| 101 | Mar | 55k | 55k |
| 102 | Jan | 60k | 60k |
| 102 | Mar | 65k | 65k |
| 103 | Jan | 40k | 40k |
| 103 | Feb | 42k | 42k |
| 103 | May | 45k | 45k |

Case 4: previous → current { sliding window }

last_value('salary').over()

Window.partitionBy('emp-id')

• orderBy('month')

• rowsBetween(-1, 0)

| emp-id | month | salary | last_sal |
|--------|-------|--------|----------|
| 101 | Jan | 50k | 50k |
| 101 | Feb | 52k | 52k |
| 101 | Mar | 55k | 55k |
| 102 | Jan | 60k | 60k |
| 102 | Mar | 65k | 65k |
| 103 | Jan | 40k | 40k |
| 103 | Feb | 42k | 42k |
| 103 | Mar | 45k | 45k |

Note:

last_value() → useless with the default frame

Real-world use cases:

- final salary
- closing balance
- end of period price
- latest status per entity

nth-values:

- $\text{nth_value}(\text{col}, n)$ returns the n th-value from the window frame.
- 'n' starts from 1 { i.e. from 1st value in frame. }
- It is frame-dependent (like first_value, last_value)
- If the n th row is not visible in the frame $\rightarrow \text{NULL}$

Case 1: Default Frame

$\text{nth_value}(\text{'salary'}, 2), \text{over} ($

$\text{Window.partitionBy('emp-id')}. \text{orderBy('month')}$

O/p:

| emp-id | month | salary | nth-value |
|--------|-------|--------|-----------|
| 101 | jan | 50K | NULL |
| 101 | feb | 52K | 52K |
| 101 | mar | 55K | 52K |
| 102 | jan | 60K | NULL |
| 102 | mar | 65K | 65K |
| 103 | jan | 40K | NULL |
| 103 | feb | 42K | 42K |
| 103 | mar | 45K | 42K |

why?

- \rightarrow Jan \rightarrow only 1 row visible. no 4th
- \rightarrow Feb $\rightarrow [Jan, Feb] \rightarrow$ end of b
- \rightarrow Mar $\rightarrow [Jan, Feb, Mar] \rightarrow$ 4th

Case 2: Full Partition Frame

$\text{nth_value}(\text{"salary"}, 2), \text{over} ($

$\text{Window.partitionBy('emp-id')}$

$\cdot \text{orderBy('month')}$

$\cdot \text{rowBetween(}$

$\text{Window.unboundedPreceding}$

$\text{Window.unboundedFollowing}$

| emp-id | month | salary | nth-value |
|--------|-------|--------|-----------|
| 101 | jan | 50K | 52K |
| 101 | feb | 52K | 52K |
| 101 | mar | 55K | 55K |
| 102 | jan | 60K | 60K |
| 102 | mar | 65K | 65K |
| 103 | jan | 40K | 42K |
| 103 | feb | 42K | 42K |
| 103 | mar | 45K | 42K |

Cases: N larger than available rows

`nth_value('salary', 3).over(`

`Window.partitionBy('emp-id').orderBy('month')`

| emp_id | month | salary | nth-value |
|--------|-------|--------|-----------|
| 101 | Jan | 50k | NULL |
| 101 | Feb | 52k | NULL |
| 101 | Mar | 55k | NULL |
| 102 | Jan | 60k | NULL |
| 102 | Mar | 65k | NULL |
| 103 | Jan | 40k | NULL |
| 103 | Feb | 42k | NULL |
| 103 | Mar | 45k | NULL |

Key rules:

- ① 'N' starts from 1
- ② Depends on frame clause
- ③ Default frame
= partial visibility
- ④ Use full frame for stable results

Real world use cases

- Second purchase amt.
- Second salary revision
- Nth event in a lifecycle
- milestone tracking.