

## Rank window Functions

{ row\_numbers()  
rank()  
dense\_rank()  
percent\_rank()  
ntile(n)  
cume\_dist()

(Typically used)

(Value) 1000000

(Result) 1-1000

(Percent) 0-100

(Ntile) 1-1000

(Cumulative) 0-1000

### row\_numbers():

assigns a unique, sequential numbers to each row within a window

Starts from 1

No duplicates

Even tied values get different numbers

→ It only cares about row position after sorting

#### row\_numbers()

\*Note

must be used with

groupBy, countWindow, over()

• orderBy

(otherwise numbering is meaningless)

sample data:

emp	dept	salary
A	IT	1000
B	IT	1500
C	IT	1500
D	HR	900
E	HR	1100

## ① Basic "row\_numbers()" (global ranking):

```
df.withColumn('rn', row_number().over(Window.orderBy('salary')))
```

o/p:

emp	dept	salary	rn
D	HR	900	1
A	IT	1000	2
E	HR	1100	3
B	IT	1500	4
C	IT	1500	5

## ② row\_numbers() with "partitionBy"

```
df.withColumn('rn', row_number().over(Window.partitionBy('dept').orderBy('salary'))))
```

o/p:

emp	dept	salary	rn
A	IT	1000	1
B	IT	1500	2
C	IT	1500	3
D	HR	900	1
E	HR	1100	2

If values are equal:  
 > row\_numbers() still increments

> Order among fees is not guaranteed

### ③ Controlling tie orders

df.withColumn(

'rn',

row\_number().over(

Window.partitionBy('dept')

• orderBy('salary', 'emp')

### ④ Top-N per group:

i.e. (Top 2 salaries per department)

top = df.withColumn(

'rn',

row\_number().over(

Window.partitionBy('dept')

• ~~orderBy('salary').desc()~~

• orderBy(col('salary').desc())

top.filter(col('rn') <= 2).show()

Why row\_number() here?

• guarantees exactly N rows

• No ambiguity with ties

"row\_number()"

Best for : {Dedup, top-N}

Deduplication :

Window.partitionBy('id').orderBy(col('updated-at').descs)

id	name	updated-at
101	A	2024-01-01
101	A	2024-01-05
101	A	2024-01-10
102	B	2024-01-03
102	B	2024-01-08

This is how the duplicates are dropped using row\_number

df = df.withColumn("rn",

row\_number().over()

Window.partitionBy('id')

.orderBy(updated-at).descs)

df.filter(col('rn') == 1)

id	updated-at	rn
101	2024-01-10	1
102	2024-01-03	1

id	updated-at	rn
101	2024-01-10	1
101	2024-01-05	2
101	2024-01-01	3
102	2024-01-08	1
103	2024-01-03	2

## Rank():

- assigns a rank to focus within a window base on ordering, & it handles ties by giving the same rank

Ranks can skip numbers

fixed values → same rank}

### Sample data

Emp	dept	salary
A	IT	1000
B	IT	1500
C	IT	1500
D	HR	900
E	HR	1100

1. Same rank

2. Gaps available

i.e. gap-ranking

Ex:

Olympic medal assigning;  
if \_\_\_\_\_ & people score  
high & same, they are given with  
Gold; 3rd will be Bronze

### 1) Basic rank():

df. with Column('dept')

rank(). over()

window.orderBy('salary')

Emp	dept	salary	rn
D	HR	900	1
A	IT	1000	2
E	HR	1100	3
B	IT	1500	4
C	IT	1500	5

### 2) rank() with partitionBy:

df. with Column('dept')

rank(). over()

window.partitionBy('dept')

.orderBy('salary')

Emp	dept	salary	rn
A	IT	1000	1
B	IT	1500	2
C	IT	1500	2
D	HR	900	1
E	HR	1100	2

### 3. Top-N with rank()

Top 2 salaries per department

$df = df.withColumn('rn',$

$rank().over($

$Window.partitionBy('dept')$

$\cdot orderBy('salary')$

$\cdot orderBy(col('salary').desc())$

$df.filter(col('rn') <= 2)$

### rank() vs row\_number()

features	row_number()	rank()
unique number	Yes	No
Tie handling	X	✓
Gaps in ranking	X	✓
Best use	Dedup, exact N	LeaderBoards

### Real-world production use case of "rank()"

#### ① Leaderboards:

- Top performances
- Top selling product
- Exam ranks

#### ② Reporting:

- Show all employees with 2nd highest salary
- Top 3 scores including tie

Ex:- Olympics medal ranking is a classic example for "rank()".

## dense\_rank():

- assigns ranks to rows within a window; gives the same rank to ties, & does not skip ranks
- NO gaps in ranking.

### ① Basic "dense\_rank()"

```
df.withColumn("d8",
    dense_rank().over(
        Window.partitionBy("dept")
            .orderBy("salary")))
```

O/P

emp	dept	salary	d8
D	HR	900	1
A	IT	1000	2
E	HR	1100	3
B	IT	1500	4
C	IT	1500	4
			5

### ② dense\_rank() with partitionBy:

```
df.withColumn("d8",
    dense_rank().over(
        Window.partitionBy("dept")
            .orderBy("salary")))
```

emp	dept	salary	d8
A	IT	1000	1
B	IT	1500	2
C	IT	1500	2
D	HR	900	1
E	HR	1100	2

### rank() vs. dense\_rank()

salary	rank	dense_rank
1000	1	1
1500	2	2
1500	2	2
1600	4	3

rank() → skips  
dense\_rank() → continues

Ex:

### 1) School exam ranking:

Two students tie for 1st

Next student is 2nd, not 3rd

This is dense\_rank()

### Real-world production use cases

#### 1) Reporting dashboards

- Top categories
- Top-selling products
- Salary bands

#### 2) Analytics

- Distinct ordering
- Bucketed ranking

### dense\_rank() vs rank()

rank()	dense_rank()
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	10
11	11
12	12
13	13
14	14
15	15
16	16
17	17
18	18
19	19
20	20
21	21
22	22
23	23
24	24
25	25
26	26
27	27
28	28
29	29
30	30
31	31
32	32
33	33
34	34
35	35
36	36
37	37
38	38
39	39
40	40
41	41
42	42
43	43
44	44
45	45
46	46
47	47
48	48
49	49
50	50
51	51
52	52
53	53
54	54
55	55
56	56
57	57
58	58
59	59
60	60
61	61
62	62
63	63
64	64
65	65
66	66
67	67
68	68
69	69
70	70
71	71
72	72
73	73
74	74
75	75
76	76
77	77
78	78
79	79
80	80
81	81
82	82
83	83
84	84
85	85
86	86
87	87
88	88
89	89
90	90
91	91
92	92
93	93
94	94
95	95
96	96
97	97
98	98
99	99
100	100

### rank() vs dense\_rank()

#### rank() vs dense\_rank()

#### rank() vs dense\_rank()

rank()	dense_rank()
1	1
2	2
3	3
4	4
5	5
6	6
7	7
8	8
9	9
10	10
11	11
12	12
13	13
14	14
15	15
16	16
17	17
18	18
19	19
20	20
21	21
22	22
23	23
24	24
25	25
26	26
27	27
28	28
29	29
30	30
31	31
32	32
33	33
34	34
35	35
36	36
37	37
38	38
39	39
40	40
41	41
42	42
43	43
44	44
45	45
46	46
47	47
48	48
49	49
50	50
51	51
52	52
53	53
54	54
55	55
56	56
57	57
58	58
59	59
60	60
61	61
62	62
63	63
64	64
65	65
66	66
67	67
68	68
69	69
70	70
71	71
72	72
73	73
74	74
75	75
76	76
77	77
78	78
79	79
80	80
81	81
82	82
83	83
84	84
85	85
86	86
87	87
88	88
89	89
90	90
91	91
92	92
93	93
94	94
95	95
96	96
97	97
98	98
99	99
100	100

### rank() vs dense\_rank()

#### rank() vs dense\_rank()

#### rank() vs dense\_rank()

#### rank() vs dense\_rank()

## Percent\_rank():

tells you the relative position of a row as a percentage within a window

output range : 0.0 → 1.0

First row : 0.0

Last row : 1.0

formula: (Interview Imp)

$$\text{percent\_rank} = (\text{rank} - 1) / (\text{total\_rows} - 1)$$

• uses rank() internally

• That's why ties behave like rank() {with gaps}

### ① Basic "percent\_rank()"

df.withColumn(

'pr',

percent\_rank().over(

Window.orderBy('salary')

Emp	Salary	Rank	Percent_rank
D	900	1	0.0
A	1000	2	0.25
E	1500	3	0.50
B	1500	4	0.75
C	1500	4	0.75

### ② percent\_rank() with partitionBy

df.withColumn(

'pr',

percent\_rank().over(

Window.partitionBy('dept')

• orderBy('salary')

Emp	Salary	Dept	Pr
A	1000	IT	0.00
B	1500	IT	0.50
C	1500	IT	0.50
D	900	HR	0.00
E	1100	HR	0.50

## percent\_rank() vs cume\_dist()

Feature	percent_rank	cume_dist
formula-based	✓	✗
first row	0.0	> 0
last row	1.0	1.0
Tie handling	rank-based	count-based
Best for	Relative position	Coverage

### Real-world use cases:

#### ① Distributed analysis

- salary percentiles
- score normalization

#### ② Threshold filtering

$$\{ \text{percent\_rank} \leq 0.1 \} \rightarrow (\text{top } 10\%)$$

## ntile(n):



- divides the rows in a window into 'n' approximately equal buckets & assigns a bucket number to each row.
  - Buckets are numbered 1 to n
  - rows are distributed as evenly as possible.
  - order matters.

split ordered data into N groups.

### data

emp	dept	salary
A	29	1000
B	29	1200
C	29	1500
D	29	1800
E	39	2000

(Total rows = 5)

How spark distributes rows?

#### Rules:

1. Rows are ordered
2. rows are divided as evenly as possible.
3. earlier buckets get extra rows.

#### ① Basic ntile(2)

df.withColumn("bucket", ntile(2).over(Window.orderBy("salary")))

emp	dept	salary	bucket
A	29	1000	1
B	29	1200	1
C	29	1500	1
D	29	1800	2
E	39	2000	2

#### Note:

first bucket may have more rows if rows aren't divisible evenly

## ② ntile(4):

df.withColumn(  
 'bucket',  
 ntile(4).over(  
 Window.orderBy('salary'))

O/P:

emp	salary	bucket
A	1000	1
B	1200	1
C	1500	2
D	1800	3
E	2000	4

## ③ ntile(n) with partitionBy:

df.withColumn(  
 'bucket',  
 ntile(2).over(  
 Window.partitionBy('dept')  
 .orderBy('salary'))

O/P:

emp	dept	salary	bucket
A	IT	1000	1
B	IT	1200	1
C	IT	1200	2
D	HR	900	1
E	HR	1100	2

→ Buckets are calculated independently per department

### key behaviours:

- Uses row positions (not values)
  - salary gaps don't matter
- orderBy is mandatory
  - without order → buckets meaningless
- Bucket sizes differ by at most 1
  - Balanced distribution

### Real-world use cases:

#### ① Customer segmentation

- top 25%, middle 50%, bottom 25%

② salary bands

· Quartiles/deciles

: (1) q1/q3

(2) median

③ Data samples

· evenly split data

Cume\_dist():

→ returns the cumulative distribution of a row within a window

{What fraction of rows have values less than or equal to the current row?

→ output range: (0,1)

last row is always 1.0

Ties are handled naturally

How cume\_dist() is calculated?

↓  
cume\_dist =  $\frac{\text{number of rows} \leq \text{current row}}{\text{total rows}}$

Sample Data:

emp	salary
D	900
A	1000
E	1100
B	1500
C	1500

Total rows = 5

df.withColumn(

'cd'

cume\_dist().over(  
Window.orderBy('salary'))

emp	salary	cd
D	900	0.20
A	1000	0.40
E	1100	0.60
B	1500	1.00
C	1500	1.00

Why both B and C are 1.00?

All 5 rows have salary  $\leq$  1500

cume-dist()

tells : up to this row,

how much of the data is already covered?

Salaries (sorted) : 900, 1000, 1100, 1500, 1500

(total rows = 5)

→ row by row, thinking:

Row 1 → salary = 900

- Rows  $\leq$  900  $\rightarrow$  1
- Total rows  $\rightarrow$  5

$$\text{cume-dist} = 1/5 = 0.20$$

Row 2 → salary = 1000

$$\begin{aligned} \cdot \text{Rows } <= 1000 &\rightarrow 2 \\ \text{cume-dist} &= 2/5 = 0.40 \end{aligned}$$

Row 3 → salary = 1100

$$\cdot \text{Rows } <= 1100 = 3$$

$$\text{cume-dist} = 3/5 = 0.60$$

Row 5 → salary = 1500

$$\text{rows } <= 1500 \Rightarrow 5$$

$$\text{cume-dist} = 5/5 = 1.00$$

Row 4 → salary = 1500

$$\cdot \text{Rows } <= 1500 = 5$$

Why 5?

Because two rows  
have salary = 1500

$$\text{cume-dist} = 5/5 = 1.00$$

## ② cume\_dist() with partitionBy

df.withColumn(

'cd',

cume\_dist().over(

Window.partitionBy('dept')

.orderBy('salary')

)

)

dp:

emp	dept	salary	cume-dist
A	IT	1000	0.33 (1/3)
B	IT	1500	1.00 (3/3)
C	IT	1500	1.00 (3/3)
D	HR	900	0.50 (1/2)
E	HR	1100	1.00 (2/2)

feature	percentileanks	cume_dist()
formula	$(rank - 1) / (n - 1)$	$\text{rows} \leftarrow \text{current} / n$
first row	0.0	> 0
last row	1.0	1.0
tie behaviour	rank-based	count-based
Best for	relative position	coverage

Real-world use cases:

① percentile thresholds  
(cume\_dist <= 0.8)  
→ bottom 80%

② Distribution analysis  
• sales distribution  
• scores coverage