

sorts() / orderBy()

- we can use either `sorts()` or `orderBy()` function of `pyspark DataFrame` by ascending / descending order based on single or multiple columns.

$\Rightarrow df.sort("col") ; df.sort("col1", "col2")$

$\Rightarrow df.orderBy("col1") ; df.orderBy("col1", "col2")$

$df.sort(df.salary.desc()) ; df.orderBy(col("col").desc)$

$df.sort(col("col2").desc()) ; df.orderBy(col("col").asc)$

GroupBy

- It allows you to perform aggregate functions on groups of rows, rather than on individual rows, enabling you to summarize data & generate aggregate statistics.

aggregate functions (agg())

`count()` `mean()` `max()` `min()`

`sum()` `avg()`

similar to groupBy in SQL

$\Rightarrow [df.groupBy("col1")]$

df.groupBy("col1").count()

df.groupBy("col1")\

• agg(

sum("col2").alias("col1"),

avg("col2").alias("col2"),

max("col3").alias("col3"),

(col1, col2) partitionBy : (col1, col2) partitionBy

df.groupBy("col1")\

• agg(

sum("col2").alias("col1"),

)\

• filter("col1")\

• orderBy("col1")\

• limit(5)

which groupBy + having

df.groupBy("col1")\

• agg(

sum("col2").alias("col1"),

)\

• filter("col1")

((col1) >= 10) &

Joins

- PySpark joins is used to combine two datasets.
- PySpark join operation combines data from two or more datasets based on a common column or key.

Innes Join,
left join,
right join,

full outer join.
left semi join
left anti join

cross join
self join

All join types in PySpark

Standard joins

- ① Innes join
- ② Left join
(left outer join)
- ③ right join
(right outer join)
- ④ Full outer join

Special joins

- ⑤ left semi join
- ⑥ left anti join

NO;
right semi X
right anti X

Advanced joins

- ⑦ Cross join
(Cartesian join)
- ⑧ self join

- ⑨ Broadcast join
(optimization technique)
(not a join type)

① Innes Join :

- It returns only the rows that have matching keys in Both tables

→ Match exists on both side : keep row
No match on either side : drop row

orders

| o_id | cust_id | amt |
|------|---------|-----|
| 1 | 101 | 500 |
| 2 | 102 | 700 |
| 3 | 103 | 300 |

customers

| cust_id | country |
|---------|---------|
| 101 | India |
| 103 | USA |
| 104 | UK |

Inner join Syntax

- ① Same column name

```
df-orders.join(df-customers, "cust_id", "inner")
```

- ② Different column names

```
df-orders.join(
```

```
df-customers,
```

```
"df-orders.cust_id" == "df-customers.id"
```

```
"inner")
```

Q1

| o_id | cust_id | amt | country |
|------|---------|-----|---------|
| 1 | 101 | 500 | India |
| 3 | 103 | 300 | USA |

cust_id = 102 ✗ Nominal
cust_id = 104 ✗ Nominal

Ans 1 : who did no item delivery
Ans 2 : who to the no. obtainable

① Left join (left outer join)

- returns ALL rows from the left table, & only matching rows from the right table.
- If there's no match on the right \rightarrow NULLs are filled.

```
orders JOIN (customers, "cust-id", "left")  
orders JOIN (customers,  
            orders.cust-id == customers.id,  
            "left")
```

o/p: ~~before join~~ = ~~before join~~

| o-id | cust-id | amt | country |
|------|---------|-----|---------|
| 1 | 101 | 500 | India |
| 2 | 102 | 700 | NULL |
| 3 | 103 | 300 | USA |

③ right join (right outer join)

- returns ALL rows from the right table, & only matching rows from the left table.
- If there's no match on left \rightarrow NULLs are filled.

```
orders JOIN (customers, "cust-id", "right")  
orders JOIN (customers,  
            orders.cust-id == customers.id,  
            "right")
```

o/p:

| cust-id | country | o-id | amt |
|---------|---------|------|------|
| 101 | India | 1 | 500 |
| 103 | USA | NULL | 300 |
| 104 | UK | NULL | NULL |

④ Full Outer Join:

- returns ALL rows from Both tables
- If a row has no match on the other side → missing columns become NULL

Orders Customers

| o-id | cust-id | amt |
|------|---------|-----|
| 1 | 101 | 500 |
| 2 | 102 | 700 |
| 3 | 103 | 300 |

| cust-id | Country |
|---------|---------|
| 101 | India |
| 103 | USA |
| 104 | UK |

Syntax:

```
orders.join(customers, "cust-id", "Outer")
```

orders.join(

customers,

orders.cust-id == customers.id,

"Outer")

dp:

| o-id | cust-id | amt | Country |
|------|---------|------|---------|
| 1 | 101 | 500 | India |
| 2 | 102 | 700 | NULL |
| 3 | 103 | 300 | USA |
| NULL | 104 | NULL | UK |

② Left semi:

- Left Semi Join returns rows from Left table that having a match in the right table.
- No columns from the right table are returned.

Syntax:

```
orders.join(customers, "cust-id", "left_semi")
```

```
orders.join(
```

customers,

orders.cust-id == customers.id,

"left_semi"

```
)
```

left-semi - o/p:

| o-id | cust-id | amt |
|------|---------|-----|
| 1 | 101 | 500 |
| 3 | 103 | 300 |

inner join - o/p:

| o-id | cust-id | amt | country |
|------|---------|-----|---------|
| 1 | 101 | 500 | India |
| 3 | 103 | 300 | USA |

{ returns only common rows that are from only left-table }

{ returns common rows of respective all columns from both tables }

similar to exists in SQL

```
Select * From orders o  
where exists (
```

```
select 1
```

```
From customers c
```

```
where o.cust-id = c.cust-id
```

⑥ Left Anti Join

- returns rows from the left table that do not have a match in right table and for it also left filter left most applicable is 01
- No columns from the right table are returned

orders

| o-id | cust-id | amt |
|------|---------|-----|
| 1 | 101 | 500 |
| 2 | 102 | 700 |
| 3 | 103 | 300 |

customers

| cust-id | country |
|---------|---------|
| 101 | India |
| 102 | USA |

Syntax:

```
orders.join(customers, 'custId', 'left-anti')
```

orders.join(

customers,

orders.cust-id == customers.id,

| left-anti | | |
|-----------|-----|-----|
| 101 | 101 | 500 |
| 102 | 102 | 700 |

| find | bi-table | 610 |
|------|----------|-----|
| 102 | 101 | 1 |
| 100 | 801 | 8 |

Q.P.

| o-id | cust-id | amt |
|------|---------|-----|
| 2 | 102 | 700 |

left semi → rows that have a match

left anti → rows that do not have a match

In SQL;

semi = exists

anti = not exists

SQL similar code:

```
↓  
Select *  
From orders o  
Where not exists (  
    Select 1  
    From customers c  
    Where o.customer_id = c.customer_id  
)
```

② Cross Join: (Cartesian join)

Returns every possible combination of rows from both tables

(left rows \times right rows = output rows)

```
orders.crossJoin(customers)
```

③ Self Join:

is when a DataFrame is joined with itself to compare rows within the same table

Note:

Same table, different roles

You treat one DataFrame as two logical tables using aliases

employees

| emp-id | emp-name | manager-id |
|--------|----------|------------|
| 1 | A | NULL |
| 2 | B | 1 |
| 3 | C | 1 |
| 4 | D | 2 |

Syntax:

$e = \text{employees.alias("e")}$

$m = \text{employees.alias("m")}$

$\text{result} = e.\text{join}($

$m,$

$e.\text{manager_id} = m.\text{emp_id},$

"left"

$)$

$\text{Select}($

$e.\text{emp_name}.alias(\text{"employee"})$

$m.\text{emp_name}.alias(\text{"manager"})$

Union:

→ Union is used to append rows of one data frame to another.

[vertical combine, not side-by-side like joins]

{ union → stack rows }
 { join → add columns }

① Union:

- combines rows from both data frames
- keeps duplicates
- Column order & data types must match

Syntax : `df1.union(df2)`

| id | name |
|----|------|
| 1 | A |
| 2 | B |

| id | name |
|----|------|
| 2 | B |
| 3 | C |

| id | name |
|----|------|
| 1 | A |
| 2 | B |
| 2 | B |
| 3 | C |

- Duplicates are kept

What if ?

`df1.union(df2)`

df1

| id | name |
|----|------|
| 1 | A |

df2

| name | id |
|------|----|
| B | 2 |

| id | name |
|----|------|
| 1 | A |
| B | 2 |

$df_1 : id, name$; $df_2 : name, id$

~~df1~~ $df_1.union(df_2)$

⇒ Data coassumption happens silently

⇒ NO errors; but values are swapped.

Why spark does Not throw errors?

- ① Column count \Rightarrow same
- ② Data types compatible
- ③ spark does not check columns names in unions

⇒ how data types compatible?

$df_1 : id : int$; $\Rightarrow (1, 'A')$.

name : string

$df_2 : name : string$ $\Rightarrow ('B', 2)$

$id : int$

$df_1.union(df_2)$

key Rule:

{ spark finds a common super type }
{ for each column position }

⇒ spark compares column by position,
not by name

df1

position 0 → id : int
position 1 → name : string

df2

position 0 → id : int
position 1 → id : int

name : string

spark compares columns by position; not by name

↓
position 0 ⇒ df1 → int / df2 → string } common super type
position 1 ⇒ df1 → name / df2 → id } "string"
spark casts int → string

position 1 ⇒ df1 → name / string | common super type
df2 → id / int | "string"

Therefore;

df1.union(df2)

position 0 → string

position 1 → string.

we can check this by

`.printSchema()`

⇒ Due to this problem with union ~~also~~.
we basically prefer the `(unionByName())`.

`union()`

composes column by position not by names

kept duplicates

problems with `(union())`

UnionAll() → (Deprecated)

unionAll()

does not exist as a separate behavior
in modern pyspark.



In pyspark Dataframe API,

union() → already behaves like SOL UnionAll

unionAll() → is deprecated/alias and should
not be used

In SQL;

union → removes duplicates after combining rows

union All → keeps duplicates

In pyspark; present version:



union() → keeps duplicates

X unionAll() → deprecated

If we want unique rows after union()

df1.union(df2).distinct()

↳ removes duplicate rows

unionByName()

- safe & correct way to union dataframes
- unions dataframes by matching column names instead of column positions
- ⇒ This avoids silent data corruption

| df1 | |
|-----|------|
| id | name |
| 1 | A |

| df2 | |
|-----|------|
| id | name |
| 2 | B |

| df1.unionByName(df2) | |
|----------------------|------|
| C | name |
| 1 | A |
| 2 | B |

⇒ allowMissingColumns = True

cf; df1: (id, name, salary)

df2: (id, name)

→ without allowMissingColumns

df1.unionByName(df2) # Error.

→ with allowMissingColumns

df1.unionByName(df2, allowMissingColumns = True)

| id | name | salary |
|----|------|--------|
| 1 | A | 50000 |
| 2 | B | NULL |

Note:

- Column names must match
 - case-sensitive
 - spaces matter

- Data types must be compatible
 - Spark will promote types if needed

(X) Column name mismatch = error
(id vs emp-id)