① For querying the SQL Syntax in pyspark use temp view

② ASC, DESC → df.Sort(df.Name asc()).show()

③ CAST → Convert data type

df.Select(df.Salary.Cast('int'))

④ Like → operator

df. filter.(df.name ('S%.'))

⑤ where → Condition's

df.filter((df.Id =='1') & (df.name.like('S%.')).show()

⑥ Distinct

df. distinct(). show()

⑦ Drop Duplicates

df. dropDuplicates(). show()

(X) from whole table

for one Column

df. dropDuplicates(['Emp_ID']). show()

⑧ union → remove duplicates from rows
Combine two Data frames

union All → Don't remove the duplicates
Combine two Data frames

Combine two Data frame

$df_3 = df_1.$ union ALL $(df_2)$

union By name $\longrightarrow$ Combine two Data frame
based on their Column name

⑨ Joins $\longrightarrow$ Types of Joins

$Df_1.$ Join $(Df_2, Df_1.dep == Df_2.ID, 'Inner').$ Show()

⑩ upper $\longrightarrow$ function

df. widthColumn ('Name', upper (df.name)). Show()

⑪ Temp View $\longrightarrow$ we Can use this with in
a Session

df. CreateOrReplaceTempView ('Employee')

$df_1$ = Spark.Sql ("Select * from Employee")

V/s

⑫ Global temp View $\longrightarrow$ we Can use this
across the Session

df. Create or Replace Global Temp View ('Global_Emp')

$df_5$ = Spark.Sql ("Select * from Global_temp.Global_Emp")

⑬ partition by $\longrightarrow$ df $\begin{matrix} 2 \\ 2 \\ 2 \end{matrix}$ $\longrightarrow$ result

Nodes

*It used for
larger data Set
to process
fast

df. write. parquet ('path', mode = 'overwrite',
partitionby = ['Dep'])

⑭ Date format $\longrightarrow$

```
df.withColumn('date_format', date_format(lit('2024-12-21'),
                                'YYYY.MM.dd')).show()
```