

SQL & PySpark Equivalent



SQL	PySpark
SELECT	df.select()
DISTINCT	df.distinct()
WHERE	df.filter()
GROUP BY	df.groupBy()
HAVING	df.having()
ORDER BY	df.orderBy()
LIMIT	df.limit()
JOIN	df.join()

<https://www.seekhobigdata.com/>

SQL	PySpark
UNION	df.union()
INTERSECT	df.intersect()
EXCEPT	df.except()
COUNT	df.groupBy().count()
SUM	df.groupBy().sum()
AVG	df.groupBy().avg()
MIN	df.groupBy().min()
MAX	df.groupBy().max()
COALESCE	df.na.fill()

SQL	PySpark
LIKE	df.filter(col.like(pattern))
IN	df.filter(col.isin([val1, val2, ...]))
BETWEEN	df.filter((col >= min_val) & (col <= max_val))
IS NULL	df.filter(col.isNull())
IS NOT NULL	df.filter(col.isNotNull())
SUBSTRING	df.selectExpr("SUBSTRING(column, start_index, length)") df.withColumn("new_column",
CONCAT	concat(col1, col2, ...)) df.withColumn("new_column",
TRIM	trim(col))

SQL	PySpark
UPPER	<code>df.withColumn("new_column", upper(col))</code>
LOWER	<code>df.withColumn("new_column", lower(col))</code>
ROUND	<code>df.withColumn("new_column", round(col, num_decimals))</code>
DATE_FORMAT	<code>df.withColumn("new_column", date_format(col, "format"))</code>
CAST	<code>col.cast("data_type")</code> <code>row_number().over(Window.partition By().orderBy())</code>
ROW_NUMBER()	

