

Spark DataFrame Operations

Explanations & Documentation Links

Let's not get lost in the Spark DataFrame jungle! Here's a handy list of the most essential operations—including the ones we've covered in the tutorial and a few extra gems to level up your skills.

Basic Operations

Operation	Explanation	How to use	Documentation link
select	Selects specific columns from a DataFrame	<code>df.select('column_name')</code> <code>df.select(['col1', 'col2'])</code>	.select()
drop	Removes one or more columns	<code>df.drop('col_name')</code> <code>df.drop('col1', 'col2')</code>	.drop()
col	Refers to a column when used in expressions.		.col()
withColumn	Adds a new column to the DataFrame	<code>df.withColumn('new_col', some_expression)</code>	.withColumn()
withColumnRenamed	Renames a column	<code>df.withColumnRenamed('new_name', 'old_name')</code>	.withColumnRenamed()
filter	Filters rows based on a condition.	<code>df.filter(col('col_name') == 'some_value')</code> <code>df.filter(col('col_name') > 10)</code>	.filter()
alias	Assigns an alias (temporary name) to a column.	<code>df.select(col('col_name').alias('new_name'))</code>	.alias()
orderBy	Sorts rows based on a column (ascending by default)	<code>df.orderBy(col('col_name').desc())</code>	.orderBy()
sort	Sorts rows based on a column (ascending by default)	<code>df.sort(col('col_name').desc())</code>	.sort()
limit	Return only a specific number of rows.	<code>df.limit(10)</code>	.limit()
dropDuplicates	Removes duplicate rows from the DataFrame (can be based on columns)	<code>df.dropDuplicates()</code> <code>#drops all duplicate rows</code> <code>df.dropDuplicates(['col_name'])</code> <code>#drops duplicates based on a single column</code>	.dropDuplicates()
fillna	Replaces null values with a specified value.	<code>df.fillna(value)</code> <code>df.fillna({'col1': 'some value'})</code>	.fillna()
distinct	Returns a new DataFrame containing distinct rows	<code>df.distinct()</code>	.distinct()

Aggregations & Grouping

Operation	Explanation	How to use	Documentation link
groupBy	Groups the DataFrame by one or more columns	<code>df.groupBy('col_name')</code>	.groupBy()
agg	Performs aggregation on grouped data		.agg()
count	Counts the number of rows	<code>df.count()</code> <i>#counts all rows</i> <code>df.groupBy('col_name').count()</code> <i>#counts the num of rows per group</i>	.count()
sum	Computes the sum of a column.	<code>df.groupBy('col_name').sum('value')</code>	.sum()
avg	Computes the average of a column.	<code>df.groupBy('col_name').avg('value')</code>	.avg()
min	Finds the minimum value of a column.	<code>df.groupBy('col_name').min('value')</code>	.min()
max	Finds the maximum value of a column.	<code>df.groupBy('col_name').max('value')</code>	.max()

Notes:

- **.groupBy()** must be followed by an **aggregation function** (e.g., `.count()`, `.sum()`, `.avg()`).
- **.agg()** allows multiple aggregations at once

Joins & Set Operations

Operation	Explanation	How to use	Documentation link
join	Performs inner, left, right, or full joins between DataFrames.	<code>df1.join(df2, 'col_name', how = 'join type')</code>	.join()
union	Combines DataFrames (must have the same schema).	<code>df1.union(df2)</code>	.union()
intersect	Gets common rows between two DataFrames.	<code>df1.intersect(df2)</code>	.intersect()
except	Get rows present in one DataFrame but not the other.	<code>df1.except(df2)</code>	.except()

Date & Time Functions

Operation	Explanation	How to use	Documentation link
to_date	Converts a string column to a date.	<pre>df.withColumn('date_col', to_date(col('string_col'), 'yyyy-MM-dd'))</pre>	.to_date()
year / month / dayofweek	Extract year, month, or day of the week from a date.	<pre>df.withColumn('year', year(col('date_col')) df.withColumn('month', month(col('date_col')) df.withColumn('day', dayofweek(col('date_col'))</pre>	.year() .month() .dayofweek()
datediff	Calculate the difference between two dates.	<pre>df.withColumn('date_diff', datediff(col('end_date'), col('start_date'))</pre>	.datediff()