

Gentle Introduction - Overview

Keyword	Explanation	Example
Spark Session	Create a spark Session	<pre>import org.apache.spark.sql.Session val spark = Session .builder .appName("Simple Application") .getOrCreate()</pre>
Range	create DF from range	<pre>val df = spark.range(100).toDF("numbers")</pre>
Where	divisible by 2	<pre>df.where("numbers % 2 = 0").show(4)</pre>
Read option	read csv	<pre>val df = spark .read .option("inferSchema", "True") .option("header", "True") .format("csv") .load("/FileStore/tables/2015_summary-ebaee.csv")</pre>
.take(nb)	return an array of the nb first lines	<pre>df.take(3)</pre>
sort desc	sort by ascending order	<pre>df.sort(desc("colName")).show()</pre>
explain	physical plan	<pre>df.sort(desc("colName")).explain</pre>
	create an SQL table to make queries	<pre>df.createOrReplaceTempView("dfTable")</pre>
select from ...	simple query	<pre>val sqlWay = spark.sql(""" SELECT DEST_COUNTRY_NAME, count(DEST_COUNTRY_NAME) FROM dfTable GROUP BY DEST_COUNTRY_NAME ORDER BY count(1) DESC LIMIT 5 """)</pre>
count groupBy	count each col entry when grouped	<pre>val dfWay = df .groupBy("colName") .count() .orderBy(desc("colName"))</pre>
max	return the max of a col	<pre>spark.sql("SELECT max(colName) from dfTable").take(1)(0) flightDF.select(max("colName")).take(1)</pre>
sum of groups	group by a col and sum of other col	<pre>df .groupBy("colNameA") .sum("colNameB") .withColumnRenamed("sum(colNameB)", "sumColNameB") .sort(desc("sumColNameB"))</pre>