

Spark SQL Joins Cheat Sheet

Spark Joins that have a SQL Equivalent

Basic Syntax:

```
df1.join(df2, df1("column1") === df2("column2"))
      .show(false)
```

Self:

Though there is no self-join type available, we can use any of the above-mentioned join types to join DataFrame to itself.

```
private static String selfJoin() {
    DataFrame empDF = spark.createDataFrame(EMPLOYEE_DF);
    DataFrame deptDF = spark.createDataFrame(DEPT_DF);

    empDF.join(empDF, empDF.col("emp_id") === empDF.col("emp_id"))
          .select("emp_id", "dept_id")
          .as("emp1")
          .join(deptDF, empDF.col("dept_id") === deptDF.col("dept_id"))
          .select("emp1.emp_id", "dept1.dept_id")
          .as("emp1")
          .join(deptDF, empDF.col("dept_id") === deptDF.col("dept_id"))
          .select("emp1.emp_id", "dept1.dept_id")
          .show(false);
}
```

Inner:

```
empDF.join(deptDF, empDF.col("dept_id") === deptDF.col("dept_id"))
      .show(false)
```

Full outer:

```
empDF.join(deptDF, empDF.col("dept_id") === deptDF.col("dept_id"),
           "left_outer")
      .show(false)
```

Left Outer:

```
empDF.join(deptDF, empDF.col("dept_id") === deptDF.col("dept_id"),
           "left_outer")
      .show(false)
```

Spark Joins that are not in SQL!

Left semi join

Returns all rows from the left DF on records match in the right dataset on join expression, records not matched on join expression are ignored from both left and right datasets.

```
empDF.join(deptDF, empDF.col("dept_id") === deptDF.col("dept_id"),
           "left_semi")
      .show(false)
```

Left anti join

leftanti join returns only columns from the left DataFrame /Dataset for non-matched records. empDF.join(deptDF, empDF.col("dept_id") === deptDF.col("dept_id"),
 "left_anti")
 .show(false)

use the CROSS JOIN syntax to allow cartesian products between these relations, or: enable implicit cartesian products by setting the configuration variable `spark.sql.crossJoin.enabled = true`:

```
private static void crossJoin() {
    DataFrame empDF = spark.createDataFrame(EMPLOYEE_DF);
    DataFrame deptDF = spark.createDataFrame(DEPT_DF);

    empDF.join(deptDF, empDF.col("dept_id") === deptDF.col("dept_id"),
               "cross")
      .show(false);
}
```

```
private static void usingCrossJoin() {
    DataFrame empDF = spark.createDataFrame(EMPLOYEE_DF);
    DataFrame deptDF = spark.createDataFrame(DEPT_DF);

    empDF.crossJoin(deptDF)
      .show();
}
```

Spark SQL Joins Cheat Sheet

Spark Joins that are not in SQL! (cont)

```
emp  DF.c  ro  ssJ  oin  (de  ptD  F).s  ho  w(f  alse)
```
