

SparkSQL

Spark SQL允许使用SQL或熟悉的DataFrame API查询Spark程序内的**结构化数据**。DataFrames和SQL提供了访问各种数据源的常用方式，包括**Hive , Avro , Parquet , ORC , JSON和JDBC**。server模式为BI工具提供行业标准的JDBC和ODBC连接。

参考

Datasets and DataFrames

RDD、DataSet、DataFrame的区别

RDD

优点:

1. 编译时类型安全：编译时就能检查出类型错误
2. 面向对象的编程风格：直接通过类名点的方式来操作数据

缺点:

1. 序列化和反序列化的性能开销：无论是集群间的通信, 还是IO操作都需要对对象的结构和数据进行序列化和反序列化.
2. GC的性能开销：频繁的创建和销毁对象, 势必会增加GC

DataFrame

DataFrame引入了schema和off-heap，解决了RDD的缺点, 但是却丢了RDD的优点。DataFrame不是类型安全的，API也不是面向对象风格的。

- schema：RDD每一行的数据，结构都是一样的。这个结构就存储在schema中。Spark通过schame就能够读懂数据，因此在通信和IO时就只需要序列化和反序列化数据，而结构的部分就可以省略了。
- off-heap：意味着JVM堆以外的内存，这些内存直接受操作系统管理（而不是JVM）。Spark能够以二进制的形式序列化数据(不包括结构)到off-heap中，当要操作数据时，就直接操作off-heap内存。由于Spark理解schema，所以知道该如何操作。

DataSet

DataSet结合了RDD和DataFrame的优点, 并带来了一个新的概念Encoder

当序列化数据时, Encoder产生字节码与off-heap进行交互, 能够达到按需访问数据的效果, 而不用反序列化整个对象。

DataSet是Spark2的新特性, 从接口和1.6中的DataFrame完全兼容。

Quick Start

SparkSession

SparkSession是DataSet和DataFrame的API切入点, 等同于SparkContext之于RDD以及SparkStreamingContext之于流处理。

```
1 import
2
3
4
5
6
7 val spark = SparkSession
8   .builder()
9   .appName("Spark SQL basic example")
10  .config("spark.some.config.option", "some-value")
11  .getOrCreate()
12
13 // For implicit conversions like converting RDDs to DataFrames
14 import spark.implicits._
```

Creating DataFrames

可以从RDD、Hive table、或者[Spark data sources](#)中创建DataFrames。

```
1 ### spark2-shell --driver-class-path /usr/share/java/mysql-connector-java.
2
```

```

3 val jdbcDF = spark.read.format("jdbc")
4 .option("url", "jdbc:mysql://172.24.33.20:63306/scm")
5 .option("driver", "com.mysql.jdbc.Driver")
6 .option("user", "scm")
7 .option("password", "scm123qaz")
8 .option("dbtable", "hosts")
9 .load()
10
11 jdbcDF.show()
12 jdbcDF.printSchema()
13 jdbcDF.select("IP_ADDRESS").show()
14
15 ### 读HDFS上的Parquet文件
16 val parquetFile = spark.read.parquet("file:///opt/users.parquet")
17 parquetFile.show()
18

```

使用DataSet

SparkSQL支持将存在的RDD装换为Datasets



```

1 case class Person(name: String, age: Long)
2
3 // Encoders are created for case classes
4 val caseClassDS = Seq(Person("Andy", 32)).toDS()
5 caseClassDS.show()
6 // +-----+-----+
7 // |name|age|
8 // +-----+-----+
9 // |Andy| 32|
10 // +-----+-----+
11
12 // Encoders for most common types are automatically provided by importing
13 val primitiveDS = Seq(1, 2, 3).toDS()
14 primitiveDS.map(_ + 1).collect() // Returns: Array(2, 3, 4)
15
16 // DataFrames can be converted to a Dataset by providing a class. Mapping
17 val path = "examples/src/main/resources/people.json"

```

```
18 val peopleDS = spark.read.json(path).as[Person]
19 peopleDS.show()
20 // +-----+-----+
21 // |  age|  name|
22 // +-----+-----+
23 // |null|Michael|
24 // | 30|  Andy|
25 // | 19| Justin|
26 // +-----+-----+
```