



Lernmodul zur Bearbeitung großer Datenmengen mit Apache Spark

SparkSession

Starten einer SparkSession namens <name> in einem Cluster mit <clustertyp>	<code>spark = SparkSession.builder.master(clustertype).appName(name).getOrCreate()</code>
Stoppen einer SparkSession	<code>spark.stop()</code>

RDD Basics

RDD aus Collection erstellen (ParallelizeCollectionRDD)	<code>spark.sparkContext.parallelize(collection)</code>
Anzahl Partitionen eines RDDs ausgeben	<code>rdd.getNumPartitions()</code>
Gesamtes RDD anzeigen	<code>rdd.collect()</code>
Die ersten n Zeilen eines RDDs anzeigen	<code>rdd.take(n)</code>
Sortieren eines RDDs	<code>rdd.sortByKey()</code>
Gruppieren eines RDDs	<code>rdd.groupByKey()</code>

DataFrame Basics

DataFrame aus Collection erstellen	<code>spark.createDataFrame(data=collection, schema=schema)</code>
DataFrame aus CSV-Datei erstellen (mit Schema und Spaltennamen)	<code>spark.read.csv(filepath, *inferSchema=True*, *header=True*)</code>
DataFrame aus Datenbanktabelle	<code>spark.read.jdbc(url, tablename)</code>
Schema erzeugen	<code>StructType([StructField(columnname, type, nullable), ...])</code>
DataFrame anzeigen	<code>df.show()</code>
Die ersten n Zeilen eines DataFrames anzeigen	<code>df.show(n)</code>
Spalte(n) eines DataFrames auswählen	<code>df.select(columnname(s)) / df.select(df.columnname, ..)</code>
Schema eines DataFrames anzeigen	<code>df.printSchema()</code>

Statistische Informationen zu einer Spalte/ mehreren Spalten anzeigen	<code>df.describe(columnname(s))</code>
Anzahl der Einträge eines DataFrames ausgeben	<code>df.count()</code>
Sortieren eines DataFrames	<code>df.sort(columnname) / df.orderBy(columnname)</code>
Gruppieren eines DataFrames	<code>df.groupBy(columnname)</code>
SQL-Abfrage auf temporäre View eines DataFrames	<code>df.createOrReplaceTempView(viewname)</code> <code>spark.sql(query)</code>
Zusammenfügen von DataFrames über einer Spalte	<code>df1.join(df2, df1.columnname == df2.columnname, *how*)</code>
Zusammenfügen von DataFrames mit gleichem Schema	<code>df1.union(df2)</code>

Analyse, Filterung und Bereinigung von DataFrames

Filtern von DataFrames	<code>df.filter(function) / df.where(function)</code>
Duplikate entfernen	<code>df.dropDuplicates()</code>
Nullwerte ersetzen	<code>df.fillna(replacement) / df.na.fill(replacement)</code>
Nullwerte entfernen	<code>df.dropna() / df.na.drop()</code>
Einzigartige Werte einer Spalte ausgeben	<code>df.select(columnname).distinct()</code>

Transformation von DataFrames

Spalten eines DataFrames verändern	<code>df.withColumn(columnname, function)</code>
Datentyp einer Spalte ändern	<code>df.columnname.cast(type)</code>
Stringwerte ersetzen	<code>df.withColumn(columnname, regexp_replace(columnname, oldvalue, newvalue)</code>
Numerische Werte ersetzen	<code>df.withColumn(columnname, replace(columnname, oldvalue, newvalue)</code>
Spalten eines DataFrames umbenennen	<code>df.withColumnRenamed(oldname, newname)</code>
Spalte eines DataFrames entfernen	<code>df.drop(columnname)</code>
Einem DataFrame eine Spalte hinzufügen	<code>df.select(columnname(s), function.alias(columnname))</code>
User Defined Function erstellen	<code>udf(lambda function)</code>