

```
1
2 1.
3 When using Spark from Python or R, you don't write explicit JVM instructions; instead,
4 you
5 write Python and R code that Spark translates into code that it then can run on the
6 executor
7 JVMs
8
9 2.
10 The SparkSession instance is the way Spark executes
11 user-defined manipulations across the cluster
12
13 val spark = SparkSession
14     .builder()
15     .appName("Spark DF example")
16     .config("spark.master", "local")
17     .getOrCreate()
18
19 3.
20 A DataFrame is the most common Structured API and simply represents a table of data with
21 rows and columns
22 it's quite easy to convert Pandas (Python) DataFrames to Spark DataFrames
23
24 4.
25 Spark will not act on transformations until we call an action
26
27 Transformations consisting of narrow dependencies ->only one partition
28 contributes to at most one output partition
29
30 A wide dependency (or wide transformation) style transformation will have input
31 partitions
32 contributing to many output partitions
33
34 With narrow transformations, Spark will
35 automatically perform an operation called pipelining, meaning that if we specify
36 multiple filters
37 on DataFrames, they'll all be performed in-memory
38
39 5. Lazy evaluation means that Spark will wait until the very last moment to execute
40 Spark compiles this plan from
41 your raw DataFrame transformations to a streamlined physical plan that will run as
42 efficiently as
43 possible across the cluster
44
45 6.
46 val df = spark.read.format("csv").option("header","true").option("inferSchema","true")
47     .option("path","home/csv").option("mode","failfast").load()
48
49 By default, when we perform a shuffle, Spark
50 outputs 200 shuffle partitions. Let's set this value to 5 to reduce the number of the
51 output
52 partitions from the shuffle
53 spark.conf.set("spark.sql.shuffle.partitions", "5")
54 flightData2015.sort("count").take(2)
55
56 7. use sql and df are same
57
58 spark.sql("SELECT max(count) from flight_data_2015").take(1)
59 flightData2015.select(max("count")).take(1)
60
61 val maxSql = spark.sql("""
62 SELECT DEST_COUNTRY_NAME, sum(count) as destination_total
63 FROM flight_data_2015
64 GROUP BY DEST_COUNTRY_NAME
65 ORDER BY sum(count) DESC
66 LIMIT 5
```

```

61  """)
62
63  flightData2015
64  .groupBy("DEST_COUNTRY_NAME")
65  .sum("count")
66  .withColumnRenamed("sum(count)", "destination_total")
67  .sort(desc("destination_total"))
68  .limit(5)
69  .show()
70
71
72  8 create df
73  val someData = Seq(
74      Row(8, "bat"),
75      Row(64, "mouse"),
76      Row(-27, "horse")
77  )
78
79  val someSchema = List(
80      StructField("number", IntegerType, true),
81      StructField("word", StringType, true)
82  )
83
84  val someDF = spark.createDataFrame(
85      spark.sparkContext.parallelize(someData),
86      StructType(someSchema)
87  )
88
89
90  2)toDF()
91  import spark.implicits._
92  val someDF = Seq(
93      (8, "bat"),
94      (64, "mouse"),
95      (-27, "horse")
96  ).toDF("number", "word")
97
98  9. Columns represent a simple type like an integer or string, a complex type like an
    array or map
99
100  10. process of execution
101  1. Write DataFrame/Dataset/SQL Code.
102  2. If valid code, Spark converts this to a Logical Plan.
103  3. Spark transforms this Logical Plan to a Physical Plan, checking for optimizations
    along
104  the way.
105  4. Spark then executes this Physical Plan (RDD manipulations) on the cluster.
106
107  11.
108  We can either let a data source
109  define the schema (called schema-on-read) or we can define it explicitly ourselves
110
111  12. StructType made up of a number of fields, StructFields
112  StructType(StructField(DEST_COUNTRY_NAME,StringType,true),
113  StructField(ORIGIN_COUNTRY_NAME,StringType,true),
114  StructField(count,LongType,true))
115
116  If the types in the data (at runtime) do not match
117  the schema, Spark will throw an error.
118
119  val myManualSchema = StructType(Array(
120  StructField("DEST_COUNTRY_NAME", StringType, true),
121  StructField("ORIGIN_COUNTRY_NAME", StringType, true),
122  StructField("count", LongType, false,
123  Metadata.fromJson("{\"hello\":\"world\"}"))
124  ))

```

```

125 val df = spark.read.format("json").schema(myManualSchema)
126   .load("/data/flight-data/json/2015-summary.json")
127
128 13,
129 There are a lot of different ways to construct and refer to columns but the two
    simplest ways are
130 by using the col or column functions.
131
132 col("someColumnName")
133 column("someColumnName")
134 scala: $"myColumn"
135
136 expr("someCol") is equivalent to col("someCol").
137 expr("someCol - 5") => col("someCol") - 5 => expr("someCol") - 5
138 expr("(((someCol + 5) * 200) - 6) < otherCol")
139
140 14 Creating Rows
141 val myRow = Row("Hello", null, 1, false)
142 //access rows
143 myRow(0) // type Any
144 myRow(0).asInstanceOf[String] // String
145 myRow.getString(0) // String
146 myRow.getInt(2) // Int
147
148 15 DataFrame Transformations
149 We can add rows or columns
150 We can remove rows or columns
151 We can transform a row into a column (or vice versa)
152 We can change the order of rows based on the values in columns
153
154 16 Creating DataFrames
155 val df = spark.read.format("json")
156   .load("/data/flight-data/json/2015-summary.json")
157   df.createOrReplaceTempView("dfTable")
158
159 We can also create DataFrames on the fly by taking a set of rows and converting them to a
160 DataFrame
161 val myManualSchema = new StructType(Array(
162   new StructField("some", StringType, true),
163   new StructField("col", StringType, true),
164   new StructField("names", LongType, false)))
165 val myRows = Seq(Row("Hello", null, 1L))
166 val myRDD = spark.sparkContext.parallelize(myRows)
167 val myDf = spark.createDataFrame(myRDD, myManualSchema)
168 myDf.show()
169
170 val myDF = Seq(("Hello", 2, 1L)).toDF("col1", "col2", "col3")
171
172 17
173 selectExpr method when you're working with expressions in strings
174
175 You can select multiple columns by using the same style of query
176 df.select("DEST_COUNTRY_NAME", "ORIGIN_COUNTRY_NAME").show(2)
177
178 18 you can refer to columns in a number of different ways
179 df.select(
180   df.col("DEST_COUNTRY_NAME"),
181   col("DEST_COUNTRY_NAME"),
182   column("DEST_COUNTRY_NAME"),
183   'DEST_COUNTRY_NAME,
184   $"DEST_COUNTRY_NAME",
185   expr("DEST_COUNTRY_NAME"))
186   .show(2)
187
188 //but don't mix, this iwll cause error
189 df.select(col("DEST_COUNTRY_NAME"), "DEST_COUNTRY_NAME")

```

```

190
191 19 rename column names
192 df.select(expr("DEST_COUNTRY_NAME AS destination")).show(2)
193
194 //The preceding operation changes the column name back to its original name.
195 df.select(expr("DEST_COUNTRY_NAME as destination").alias("DEST_COUNTRY_NAME"))
196 .show(2)
197
198 20. Because select followed by a series of expr is such a common pattern, Spark has a
    shorthand
199 for doing this efficiently: selectExpr
200
201 df.selectExpr("DEST_COUNTRY_NAME as newColumnName", "DEST_COUNTRY_NAME").show(2)
202
203 adds a new column withinCountry to our DataFrame that specifies whether the destination
    and
204 origin are the same
205 df.selectExpr(
206     "*", // include all original columns
207     "(DEST_COUNTRY_NAME = ORIGIN_COUNTRY_NAME) as withinCountry")
208 .show(2)
209
210 21
211 we can also specify aggregations over the entire DataFrame by taking
212 advantage
213 df.selectExpr("avg(count)", "count(distinct(DEST_COUNTRY_NAME))").show(2)
214
215 22.
216 pass explicit values into Spark that are just a value (rather than a new
217 column)
218 df.select(expr("*"), lit(1).as("One")).show(2)
219
220 23. Adding Columns
221 //add numberOne
222 df.withColumn("numberOne", lit(1)).show(2)
223
224 //or by expr
225 df.withColumn("withinCountry", expr("ORIGIN_COUNTRY_NAME == DEST_COUNTRY_NAME"))
226 .show(2)
227 df.withColumn("Destination", expr("DEST_COUNTRY_NAME")).columns
228
229 24. Renaming Columns
230 This will rename the column with the name of the string in
231 the first argument to the string in the second argument:
232 df.withColumnRenamed("DEST_COUNTRY_NAME", "dest").columns
233
234 25. characters like spaces or dashes in column names
235 // column name is 'this long column-name'
236 val dfWithLongColName = df.withColumn(
237     "This Long Column-Name",
238     expr("ORIGIN_COUNTRY_NAME"))
239
240 //In this example, however, we need to use backticks because we're
241 referencing a column in an expression:
242 dfWithLongColName.selectExpr(
243     "`This Long Column-Name`",
244     "`This Long Column-Name` as `new col`")
245 .show(2)
246
247 27 Case Sensitivity
248 By default Spark is case insensitive
249 you can make Spark case sensitive by setting: set spark.sql.caseSensitive true
250
251 28
252 Removing Columns
253 df.drop("ORIGIN_COUNTRY_NAME").columns

```

```

254
255 29 Changing a Column's Type (cast)
256 convert our count column from
257 an integer to a type Long: df.withColumn("count2", col("count").cast("long"))
258
259 30.filters
260 df.filter(col("count") < 2).show(2)
261 df.where("count < 2").show(2)
262
263 Spark automatically performs all filtering operations at
264 the same time regardless of the filter ordering
265 This means that if you want to specify multiple
266 AND filters, just chain them sequentially and let Spark handle the rest:
267
268 df.where(col("count") < 2).where(col("ORIGIN_COUNTRY_NAME") != "Croatia")
269 .show(2)
270
271 31Getting Unique Rows
272 df.select("ORIGIN_COUNTRY_NAME", "DEST_COUNTRY_NAME").distinct().count()
273
274 32 sample
275 总数 256.
276
277 val seed = 5
278 val withReplacement = false
279 val fraction = 0.5
280 df.sample(withReplacement, fraction, seed).count()
281 output of 126.
282
283 33 Random Splits
284 //we'll split our DataFrame into two different
285 DataFrames by setting the weights by which we will split the DataFrame
286 // in Scala
287 val dataFrames = df.randomSplit(Array(0.25, 0.75), seed)
288 dataFrames(0).count() > dataFrames(1).count() // False
289
290 34 union
291 To append to a DataFrame, you must
292 union the original DataFrame along with the new DataFrame.
293 To union two DataFrames, you must be sure that they have the same schema and
294 number of columns; otherwise, the union will fail.
295 df.union(newDF)
296 .where("count = 1")
297 .where($"ORIGIN_COUNTRY_NAME" != "United States")
298 .show() // get all of them and we'll see our new rows at the end
299
300 you must use the != operator so that you don't just compare the unevaluated column
301
302 35
303 Sorting Rows
304
305 There are two equivalent operations to do this sort
306 and orderBy that work the exact same way
307 df.sort("count").show(5)
308 df.orderBy("count", "DEST_COUNTRY_NAME").show(5)
309
310 #use the asc and desc functions
311 df.orderBy(expr("count desc")).show(2)
312 df.orderBy(desc("count"), asc("DEST_COUNTRY_NAME")).show(2)
313 #
314
315 An advanced tip is to use asc_nulls_first, desc_nulls_first, asc_nulls_last, or
316 desc_nulls_last to specify where you would like your null values to appear
317
318 36 repartition
319 it can be worth

```

```

320 repartitioning based on that column
321
322 df.repartition(col("DEST_COUNTRY_NAME"))
323 df.repartition(5, col("DEST_COUNTRY_NAME"))
324
325 Coalesce, on the other hand, will not incur a full shuffle and will try to combine
partitions
326 df.repartition(5, col("DEST_COUNTRY_NAME")).coalesce(2)
327
328 37
329 collect some of your data to the driver in order to manipulate
330 collectDF.collect()
331 The method toLocalIterator collects partitions to the driver as an iterator
332 collectDF.toLocalIterator()
333
334 38
335 the lit function. This function converts a
336 type in another language to its corresponding Spark representation
337
338 df.select(lit(5), lit("five"), lit(5.0))
339 #
340 39 Working with Booleans
341 df.where(col("InvoiceNo").equalTo(536365))
342 .select("InvoiceNo", "Description")
343 .show(5, false)
344
345 df.where(col("InvoiceNo") === 536365)
346 .select("InvoiceNo", "Description")
347 .show(5, false)
348
349 df.where("InvoiceNo = 536365")
350 .show(5, false)
351
352 val priceFilter = col("UnitPrice") > 600
353 val descripFilter = col("Description").contains("POSTAGE")
354 df.where(col("StockCode").isin("DOT")).where(priceFilter.or(descripFilter))
355 .show()
356
357 val DOTCodeFilter = col("StockCode") === "DOT"
358 val priceFilter = col("UnitPrice") > 600
359 val descripFilter = col("Description").contains("POSTAGE")
360 df.withColumn("isExpensive", DOTCodeFilter.and(priceFilter.or(descripFilter)))
361 .where("isExpensive")
362 .select("unitPrice", "isExpensive").show(5)
363
364 df.withColumn("isExpensive", expr("NOT UnitPrice <= 250"))
365 .filter("isExpensive")
366 .select("Description", "UnitPrice").show(5)
367
368 df.withColumn("isExpensive", not(col("UnitPrice").leq(250)))
369 .filter("isExpensive")
370 .select("Description", "UnitPrice").show(5)
371
372 39 Working with Numbers
373 val fabricatedQuantity = pow(col("Quantity") * col("UnitPrice"), 2) + 5
374 df.select(expr("CustomerId"), fabricatedQuantity.alias("realQuantity")).show(2)
375
376 => df.selectExpr(
377 "CustomerId",
378 "(POWER((Quantity * UnitPrice), 2.0) + 5) as realQuantity").show(2)
379
380 we round to one decimal place:
381 df.select(round(col("UnitPrice"), 1).alias("rounded"), col("UnitPrice")).show(5)
382
383 the round function rounds up if you're exactly in between two numbers. You can
384 round down by using the bround

```

```

385 df.select(round(lit("2.5")), bround(lit("2.5"))).show(2)
386
387 #compute the correlation of two columns
388 df.stat.corr("Quantity", "UnitPrice")
389 #compute summary statistics for a column or set of columns
390 df.describe().show()
391
392 add a unique ID to each row by using the function
393 monotonically_increasing_id.
394 df.select(monotonically_increasing_id()).show(2)
395
396
397 40. Working with Strings
398 #The initcap function will
399 capitalize every word in a given string
400 df.select(initcap(col("Description"))).show(2, false)
401
402 df.select(col("Description"),
403 lower(col("Description")),
404 upper(lower(col("Description")))).show(2)
405
406 #d
407 so that length of this column should be 4 characters.
408 If length is less than 4 characters, then add 0's
409 import org.apache.spark.sql.functions.{lit, ltrim, rtrim, rpad, lpad, trim}
410 df.select(
411 ltrim(lit(" HELLO ")).as("ltrim"),
412 rtrim(lit(" HELLO ")).as("rtrim"),
413 trim(lit(" HELLO ")).as("trim"),
414 lpad(lit("HELLO"), 3, " ").as("lp"),
415 rpad(lit("HELLO"), 10, " ").as("rp")).show(2)
416
417
418 41
419
420 DESC is "WHITE HANGING HEA...|"
421 1)
422 regular replace()
423 regex_string = "BLACK|WHITE|RED|GREEN|BLUE"
424 df.select(
425 regexp_replace(col("Description"), regex_string, "COLOR").alias("color_clean"),
426 col("Description")).show(2)
427
428
429 +-----+-----+
430 | color_clean| Description|
431 +-----+-----+
432 |COLOR HANGING HEA...|WHITE HANGING HEA...|
433 | COLOR METAL LANTERN| WHITE METAL LANTERN|
434 +-----+-----+
435 2)
436 Spark also provides the translate function to replace these
437 values. This is done at the character level and will replace all
438 instances of a character with the
439 indexed character in the replacement string
440
441 val simpleColors = Seq("black", "white", "red", "green", "blue")
442 val regexString = simpleColors.map(_.toUpperCase).mkString("|")
443 // the
444 // 'L' replaced to 1, 'E' replaced to 3.
445 df.select(translate(col("Description"), "LEET", "1337"), col("Description"))
446 .show(2)
447
448
449 3)
450 pulling out the first mentioned color

```

```

451
452 # get the 'WHITE'
453 extract_str = "(BLACK|WHITE|RED|GREEN|BLUE)"
454 df.select(
455   regexp_extract(col("Description"), regexString, 1).alias("color_clean"),
456   col("Description")).show(2)
457 +-----+-----+
458 | color_clean | Description |
459 +-----+-----+
460 | WHITE      | WHITE HANGING HEA... |
461 | WHITE      | WHITE METAL LANTERN |
462 +-----+-----+
463
464 4) check for their existence
465 val containsBlack = col("Description").contains("BLACK")
466 val containsWhite = col("DESCRIPTION").contains("WHITE")
467 df.withColumn("hasSimpleColor", containsBlack.or(containsWhite))
468   .where("hasSimpleColor")
469   .select("Description").show(3, false)
470
471
472 42 Date and timestamp
473
474 dateDF = spark.range(10)\
475   .withColumn("today", current_date())\
476   .withColumn("now", current_timestamp())
477 dateDF.createOrReplaceTempView("dateTable")
478
479
480 #add , subtract dateDF
481 dateDF.select(date_sub(col("today"), 5), date_add(col("today"), 5)).show(1)
482
483 #datediff function that will return the number of days in between two dates
484 dateDF.withColumn("week_ago", date_sub(col("today"), 7))
485   .select(datediff(col("week_ago"), col("today"))).show(1)
486
487 dateDF.select(
488   to_date(lit("2016-01-01")).alias("start"),
489   to_date(lit("2017-05-22")).alias("end"))
490   .select(months_between(col("start"), col("end"))).show(1)
491
492 #The to_date function allows
493 you to convert a string to a date
494 spark.range(5).withColumn("date", lit("2017-01-01"))
495   .select(to_date(col("date"))).show(1)
496
497 #use format
498 val dateFormat = "yyyy-dd-MM"
499 val cleanDateDF = spark.range(1).select(
500   to_date(lit("2017-12-11"), dateFormat).alias("date"),
501   to_date(lit("2017-20-12"), dateFormat).alias("date2"))
502 cleanDateDF.createOrReplaceTempView("dateTable2")
503
504 #timestamp
505 dateFormat = "yyyy-dd-MM"
506 to_timestamp, which always requires a format to be specified
507 //use column date defined in previous example,
508 cleanDateDF.select(to_timestamp(col("date"), dateFormat)).show()
509
510 cleanDateDF.filter(col("date2") > lit("2017-12-12")).show()
511
512 Handle null data.
513 #Coalesce
514 first non-null value from a set of columns by
515 using the coalesce function
516 df.select(coalesce(col("Description"), col("CustomerId"))).show()

```



```

517
518 43.
519 #
520 ifnull allows you to select the second value if the first is null
521
522 nullif, which returns null if the two values are equal
523 or else returns the second if they are
524 not
525
526 nvl returns the second value if the first is null, but defaults to the first.
527
528 nvl2 returns
529 the second value if the first is not null; otherwise, it will return the last specified
    value
530
531 SELECT
532   ifnull(null, 'return_value'),
533   nullif('value', 'value'),
534   nvl(null, 'return_value'),
535   nvl2('not_null', 'return_value', "else_value")
536 FROM dfTable LIMIT 1
537
538 +-----+-----+-----+-----+
539 | a| b| c| d|
540 +-----+-----+-----+-----+
541 |return_value|null|return_value|return_value|
542 +-----+-----+-----+-----+
543
544 Using "all" drops the
545 row only if all values are null or NaN for that row
546 df.na.drop("all")
547
548 The default is to drop anyrow in which any value is null:
549 df.na.drop()
550 df.na.drop("any")
551
552 We can also apply this to certain sets of columns by passing in an array of columns:
553 //drop, 'stockCode' 和invoiceNo 都是null的列
554 df.na.drop("all", Seq("StockCode", "InvoiceNo"))
555
556 fill all null values in columns of type String
557 df.na.fill("All Null values become this string")
558
559 //把列类型Int的为空的, fill 5
560 or df.na.fill(5:Integer)
561
562 df.na.fill(5, Seq("StockCode", "InvoiceNo"))
563
564 We can also do this with with a Scala Map
565 where the key is the column name and the value is the
566 value we would like to use to fill null values:
567 // in Scala
568 val fillColValues = Map("StockCode" -> 5, "Description" -> "No Value")
569 df.na.fill(fillColValues)
570
571 replace all values in a certain column according to their current value
572 df.na.replace("Description", Map("" -> "UNKNOWN"))
573
574 44. Ordering
575 asc_nulls_first, desc_nulls_first,
576 asc_nulls_last, or desc_nulls_last to specify where you would like your null values to
577 appear in an ordered DataFrame
578
579 45
580 Working with Complex Types
581

```

```

582 1) You can think of 'structs as DataFrames within DataFrames
583 //把两个列合成一个struct,
584 val complexDF = df.select(struct("Description", "InvoiceNo").alias("complex"))
585 complexDF.createOrReplaceTempView("complexDF")
586 complexDF.select("complex.Description")
587 We can also query all values in the struct by using *.
588 complexDF.select("complex.*")
589
590 2) array
591 df.select(split(col("Description"), " ")).show(2)
592 +-----+
593 |split(Description, )|
594 +-----+
595 | [WHITE, HANGING, ...|
596 | [WHITE, METAL, LA...|
597 +-----+
598
599
600 Spark allows us to manipulate this complex type as another
601 column
602
603 df.select(split(col("Description"), " ").alias("array_col"))
604 .selectExpr("array_col[0]").show(2)
605
606 This gives us the following result:
607 +-----+
608 |array_col[0]|
609 +-----+
610 | WHITE|
611 | WHITE|
612 +-----+
613 Array Length
614 //use size() get length
615 df.select(size(split(col("Description"), " "))).show(2) // shows 5 and 3
616
617 # array contains
618 df.select(array_contains(split(col("Description"), " "), "WHITE")).show(2) => true
619
620 The explode function takes a column that consists of arrays and creates one row (with
621 the rest of
622 the values duplicated) per value in the array
623
624 #use array in desc, explode two multiple rows.
625 df.withColumn("splitted", split(col("Description"), " "))
626 .withColumn("exploded", explode(col("splitted")))
627 .select("Description", "InvoiceNo", "exploded").show(2)
628
629 [WHITE, HANGING] explode后变成两列
630 +-----+-----+-----+
631 | Description|InvoiceNo|exploded|
632 +-----+-----+-----+
633 |WHITE HANGING HEA...| 536365| WHITE|
634 |WHITE HANGING HEA...| 536365| HANGING|
635 +-----+-----+-----+
636
637 Maps
638 // in Scala
639 import org.apache.spark.sql.functions.map
640 df.select(map(col("Description"), col("InvoiceNo")).alias("complex_map")).show(2)
641 +-----+
642 | complex_map|
643 +-----+
644 |Map(WHITE HANGING...|
645 |Map(WHITE METAL L...|
646 +-----+
647 //query by key 得到invoiceNo

```

```

647 df.select(map(col("Description"), col("InvoiceNo")).alias("complex_map"))
648 .selectExpr("complex_map['WHITE METAL LANTERN']").show(2)
649 +-----+
650 |complex_map[WHITE METAL LANTERN]|
651 +-----+
652 | null|
653 | 536365|
654 +-----+
655 You can also explode map types, which will turn them into columns:
656 df.select(map(col("Description"), col("InvoiceNo")).alias("complex_map"))
657 .selectExpr("explode(complex_map)").show(2)
658 +-----+-----+
659 | key| value|
660 +-----+-----+
661 |WHITE HANGING HEA...|536365|
662 | WHITE METAL LANTERN|536365|
663 +-----+-----+
664
665 46 JSON process
666
667 get_json_object to inline query a JSON object
668 json_tuple if this object has only one level of nesting:
669
670 val jsonDF = spark.range(1).selectExpr("""
671 '{"myJSONKey" : {"myJSONValue" : [1, 2, 3]}}' as jsonString""")
672 //get the array inside json
673 jsonDF.select(
674 get_json_object(col("jsonString"), "$.myJSONKey.myJSONValue[1]") as "column",
675 json_tuple(col("jsonString"), "myJSONKey")).show(2)
676
677 This results in the following table:
678 +-----+
679 |column| c0|
680 +-----+
681 | 2|{"myJSONValue":[1...|
682 +-----+
683
684 turn a StructType into a JSON string by using the to_json
685 df.selectExpr("(InvoiceNo, Description) as myStruct")
686 .select(to_json(col("myStruct")))
687 #
688
689 from_json function to parse this (or other JSON data) back in
690 val parseSchema = new StructType(Array(
691 new StructField("InvoiceNo",StringType,true),
692 new StructField("Description",StringType,true)))
693 df.selectExpr("(InvoiceNo, Description) as myStruct")
694 .select(to_json(col("myStruct")).alias("newJSON"))
695 .select(from_json(col("newJSON"), parseSchema), col("newJSON")).show(2)
696
697 47 UFD
698 Spark will serialize the function on the
699 driver and transfer it over the network to all executor processes
700
701 If the function is written in Python, something quite different happens. Spark starts a
702 Python
703 process on the worker, serializes all of the data to a format that Python can understand
704 (remember, it was in the JVM earlier), executes the function row by row on that data in
705 the
706 Python process, and then finally returns the results of the row operations to the JVM
707 and Spark.
708
709 #register
710 val power3udf = udf(power3(_:Double):Double)
711 udfExampleDF.select(power3udf(col("num"))).show()
712

```

```

710 #another wasy
711
712 spark.udf.register("power3", power3(_:Double):Double)
713 udfExampleDF.selectExpr("power3(num)").show(2)
714
715 return type in the following function to be a DoubleType:
716 spark.udf.register("power3py", power3, DoubleType())
717
718 48 Join
719
720 1) count
721 df.select(count("StockCode")).show() // 541909
722
723 df.select(countDistinct("StockCode")).show() // 4070
724 df.select(approx_count_distinct("StockCode", 0.1)).show() // 3364
725
726 first and last
727 df.select(first("StockCode"), last("StockCode")).show()
728
729 df.select(min("Quantity"), max("Quantity")).show()
730
731 df.select(sum("Quantity")).show() // 5176450
732
733 df.select(sumDistinct("Quantity")).show() // 29310
734
735 df.select(
736   count("Quantity").alias("total_transactions"),
737   sum("Quantity").alias("total_purchases"),
738   avg("Quantity").alias("avg_purchases"),
739   expr("mean(Quantity)").alias("mean_purchases"))
740
741 The variance is the average
742 of the squared differences from the mean, and the standard deviation is the square root
of the
743 variance
744 df.select(var_pop("Quantity"), var_samp("Quantity"),
745   stddev_pop("Quantity"), stddev_samp("Quantity")).show()
746
747 Skewness
748 measures the asymmetry of the values in your data around the mean, whereas kurtosis is a
749 measure of the tail of data
750 df.select(skewness("Quantity"), kurtosis("Quantity")).show()
751
752
753 df.select(corr("InvoiceNo", "Quantity"), covar_samp("InvoiceNo", "Quantity"),
754   covar_pop("InvoiceNo", "Quantity")).show()
755 #
756 Aggregating to Complex Types
757 import org.apache.spark.sql.functions.{collect_set, collect_list}
758 df.agg(collect_set("Country"), collect_list("Country")).show()
759
760
761 df.groupBy("InvoiceNo", "CustomerId").count().show()
762 +-----+-----+-----+
763 |InvoiceNo|CustomerId|count|
764 +-----+-----+-----+
765 | 536846| 14573| 76|
766 ...
767 | C544318| 12989| 1|
768 +-----+-----+-----+
769 # Grouping with Expressions
770 df.groupBy("InvoiceNo").agg(
771   count("Quantity").alias("quan"),
772   expr("count(Quantity)").show()
773   |InvoiceNo|quan|count(Quantity)|
774 +-----+-----+-----+

```

```

775 | 536596| 6| 6|
776 ...
777 | C542604| 8| 8|
778 +-----+-----+-----+
779
780
781 Grouping with Maps
782 df.groupBy("InvoiceNo").agg("Quantity"->"avg", "Quantity"->"stddev_pop").show()
783 +-----+-----+-----+
784 |InvoiceNo| avg(Quantity)|stddev_pop(Quantity)|
785 +-----+-----+-----+
786 | 536596| 1.5| 1.1180339887498947|
787 ...
788 | C542604| -8.0| 15.173990905493518|
789 +-----+-----+-----+
790 #partitionby is unrelated to the partitioning scheme concept that we have covered thus
791 far.
792 It's just a
793 similar concept that describes how we will be breaking up our group
794 val windowSpec = Window
795 .partitionBy("CustomerId", "date")
796 .orderBy(col("Quantity").desc).rowsBetween(Window.unboundedPreceding, Window.currentRow)
797 #establishing the maximum purchase quantity over all time.
798 maxPurchaseQuantity = max(col("Quantity")).over(windowSpec)
799
800 val purchaseDenseRank = dense_rank().over(windowSpec)
801
802 Now we can perform a select to
803 view the calculated window values
804
805 dfWithDate.where("CustomerId IS NOT NULL").orderBy("CustomerId")
806 .select(
807   col("CustomerId"),
808   col("date"),
809   col("Quantity"),
810   purchaseRank.alias("quantityRank"),
811   purchaseDenseRank.alias("quantityDenseRank"),
812   maxPurchaseQuantity.alias("maxPurchaseQuantity")).show()
813
814
815 you also want to include the total number of items, regardless of
816 customer or stock code? With a conventional group-by statement, this would be
817 impossible. But,
818 it's simple with grouping sets: we simply specify that we would like to aggregate at
819 that level
820
821 SELECT CustomerId, stockCode, sum(Quantity) FROM dfNoNull
822 GROUP BY customerId, stockCode GROUPING SETS((customerId, stockCode),())
823 ORDER BY CustomerId DESC, stockCode DESC
824
825 49 Rollup
826
827 creates a new DataFrame that includes the grand total over all dates, the
828 grand total for each 'date' in the DataFrame, and the subtotal for each 'country' on
829 each date in the
830 DataFrame:
831
832 val rolledUpDF = dfNoNull.rollup("Date", "Country").agg(sum("Quantity"))
833 .selectExpr("Date", "Country", "`sum(Quantity)` as total_quantity")
834 .orderBy("Date")
835 rolledUpDF.show()
836
837 +-----+-----+-----+
838 | Date| Country|total_quantity|
839 +-----+-----+-----+

```

```

837 | null| null| 5176450|
838 |2010-12-01|United Kingdom| 23949|
839 |2010-12-01| Germany| 117|
840 |2010-12-01| France| 449|
841 ...
842 |2010-12-03| France| 239|
843 |2010-12-03| Italy| 164|
844 |2010-12-03| Belgium| 528|
845 +-----+-----+-----+
846
847 Cube
848 across all dimensions
849 The total across all dates and countries
850 The total for each date across all countries
851 The total for each country on each date
852 The total for each country across all dates
853 dfNonNull.cube("Date", "Country").agg(sum(col("Quantity")))
854 .select("Date", "Country", "sum(Quantity)").orderBy("Date").show()
855
856 grouping_id,
857 which gives us a column specifying the level of aggregation that we have in our result
858 set
859 dfNonNull.cube("customerId", "stockCode").agg(grouping_id(), sum("Quantity"))
860 .orderBy(expr("grouping_id()).desc)
861 .show()
862
863 51.
864 User-Defined Aggregation Functions
865 52. join
866
867 val joinExpression = person.col("graduate_program") === graduateProgram.col("id")
868 person.join(graduateProgram, joinExpression).show()
869
870 or
871 person.join(graduateProgram, joinExpression, joinType).show()
872
873 53 Outer Joins
874 If there is no equivalent row in either the left or
875 right DataFrame, Spark will insert null:
876
877 joinType = "outer"
878 person.join(graduateProgram, joinExpression, joinType).show()
879
880 54
881 Left Outer Joins
882 includes all rows from
883 the left DataFrame If there is no equivalent row in the right DataFrame, Spark will
884 insert null
885
886 joinType = "left_outer"
887 graduateProgram.join(person, joinExpression, joinType).show()
888
889 SELECT * FROM graduateProgram LEFT OUTER JOIN person
890 ON person.graduate_program = graduateProgram.id
891
892 55.
893 joinType = "right_outer"
894 person.join(graduateProgram, joinExpression, joinType).show()
895
896 56 left-semi
897 They do not actually include any values
898 from the right DataFrame. They only compare values to see if the value exists in the
899 second
900 DataFrame. If the value does exist, those rows will be kept in the result, even if
901 there are

```

```

899 duplicate keys in the left DataFrame
900 joinType = "left_semi"
901 graduateProgram.join(person, joinExpression, joinType).show()
902 57 Left Anti Joins
903 joinType = "left_anti"
904
905 rather than keeping the values that exist in the second
906 DataFrame, they keep only the values that do not have a corresponding key in the second
907 DataFrame
908
909 graduateProgram.join(person, joinExpression, joinType).show()
910
911 58
912 person.crossJoin(graduateProgram).show()
913 joinType = "cross"
914 graduateProgram.join(person, joinExpression, joinType).show()
915 59
916 person.join(gradProgramDupe, joinExpr).select("graduate_program").show()
917 Given the previous code snippet, we will receive an error. In this particular example,
918 Spark
919 generates this message:
920 org.apache.spark.sql.AnalysisException: Reference 'graduate_program' is
921 ambiguous, could be: graduate_program#40, graduate_program#1079.;
922
923 Joins on Complex Types
924 person.withColumnRenamed("id", "personId")
925 .join(sparkStatus, expr("array_contains(spark_status, id)").show()
926
927 Approach 1: Different join expression
928 When you have two keys that have the same name, probably the easiest fix is to change
929 the join
930 expression from a Boolean expression to a string or sequence
931
932 Approach 2: Dropping the column after the join
933 person.join(gradProgramDupe, joinExpr).drop(person.col("graduate_program"))
934 .select("graduate_program").show()
935
936 Approach 3: Renaming a column before the join
937 val gradProgram3 = graduateProgram.withColumnRenamed("id", "grad_id")
938 val joinExpr = person.col("graduate_program") === gradProgram3.col("grad_id")
939 person.join(gradProgram3, joinExpr).show()
940
941 Big table-to-small table
942 we can use a big table-to-big
943 table communication strategy, it can often be more efficient to use a broadcast join
944
945 With the DataFrame API, we can also explicitly give the optimizer a hint that we would
946 like to
947 use a broadcast join by using the correct function around the small DataFrame in
948 question.
949
950 import org.apache.spark.sql.functions.broadcast
951 val joinExpr = person.col("graduate_program") === graduateProgram.col("id")
952 person.join(broadcast(graduateProgram), joinExpr).explain()
953
954 val dbDataFrame = spark.read.format("jdbc")
955 .option("url", url).option("dbtable", tablename).option("driver", driver)
956 .option("numPartitions", 10).load()
957
958 Partitioning based on a sliding window
959 // in Scala
960 val colName = "count"
961 val lowerBound = 0L
962 val upperBound = 348113L // this is the max count in our database

```

```

961 val numPartitions = 10
962 spark.read.jdbc(url,tablename,colName,lowerBound,upperBound,numPartitions,props)
963 .count() // 255
964
965
966 spark.read.textFile("/data/flight-data/csv/2010-summary.csv")
967 .selectExpr("split(value, ',') as rows").show()
968
969 csvFile.limit(10).select("DEST_COUNTRY_NAME", "count")
970 .write.partitionBy("count").text("/tmp/five-csv-files2.csv")
971
972 #write with partition
973 csvFile.limit(10).write.mode("overwrite").partitionBy("DEST_COUNTRY_NAME")
974 .save("/tmp/partitioned-files.parquet")
975 #
976
977 Bucketing is another file organization approach with which you can control the data
978 that is
979 specifically written to each file. This can help avoid shuffles later when you go to
980 read the data
981 because data with the same bucket ID will all be grouped together into one physical
982 partition
983
984 CSV files do not support complex types, whereas Parquet and ORC do.
985
986 df.write.option("maxRecordsPerFile", 5000), Spark
987 will ensure that files will contain at most 5,000 records.
988
989 SparkSQL Thrift JDBC/ODBC Server
990 Spark provides a Java Database Connectivity (JDBC) interface by which either you or a
991 remote
992 program connects to the Spark driver in order to execute Spark SQL queries
993 ./sbin/start-thriftserver.sh
994
995 You can then test this connection by running the following commands:
996 beeline> !connect jdbc:hive2://localhost:10000
997
998 When you define a table from files on disk, you are defining an unmanaged table
999
1000 When
1001 you use saveAsTable on a DataFrame, you are creating a managed table
1002
1003 you can also see tables in a specific
1004 database by using the query show tables IN databaseName
1005
1006 Creating Tables
1007 CREATE TABLE flights (
1008   DEST_COUNTRY_NAME STRING, ORIGIN_COUNTRY_NAME STRING, count LONG)
1009 USING JSON OPTIONS (path '/data/flight-data/json/2015-summary.json')
1010
1011 It is possible to create a table from a query as well:
1012 CREATE TABLE flights_from_select USING parquet AS SELECT * FROM flights
1013
1014 we create an unmanaged table. Spark will
1015 manage the table's metadata; however, the files are not managed by Spark at all. You
1016 create this
1017 table by using the CREATE EXTERNAL TABLE statement
1018 CREATE EXTERNAL TABLE hive_flights (
1019   DEST_COUNTRY_NAME STRING, ORIGIN_COUNTRY_NAME STRING, count LONG)
1020 ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' LOCATION '/data/flight-data-hive/'
1021
1022 INSERT INTO partitioned_flights
1023 PARTITION (DEST_COUNTRY_NAME="UNITED STATES")

```



```

1022 SELECT count, ORIGIN_COUNTRY_NAME FROM flights
1023 WHERE DEST_COUNTRY_NAME='UNITED STATES' LIMIT 12
1024
1025 You can also create an external table from a select clause:
1026 CREATE EXTERNAL TABLE hive_flights_2
1027 ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
1028 LOCATION '/data/flight-data-hive/' AS SELECT * FROM flights
1029
1030 65. Describing Table Metadata
1031
1032 DESCRIBE TABLE flights_csv
1033 SHOW PARTITIONS partitioned_flights
1034
1035 66.Refreshing Table Metadata
1036 REFRESH table partitioned_flights
1037 MSCK REPAIR TABLE partitioned_flights
1038
1039 67
1040 DROP TABLE IF EXISTS flights_csv;
1041
1042 Dropping unmanaged tables
1043 If you are dropping an unmanaged table (e.g., hive_flights), no data will be removed
1044 but you
1045 will no longer be able to refer to this data by the table name
1046
1047 68. View
1048 A view specifies a set
1049 of transformations on top of an existing table—basically just saved query plans
1050 CREATE VIEW just_usa_view AS
1051 SELECT * FROM flights WHERE dest_country_name = 'United States'
1052
1053 you can create temporary views that are available only during the current session
1054 CREATE TEMP VIEW just_usa_view_temp AS
1055 SELECT * FROM flights WHERE dest_country_name = 'United States'
1056
1057 or CREATE OR REPLACE TEMP VIEW
1058
1059 views are equivalent to creating a new DataFrame from an existing DataFrame
1060 val flights = spark.read.format("json")
1061 .load("/data/flight-data/json/2015-summary.json")
1062 val just_usa_df = flights.where("dest_country_name = 'United States'")
1063 just_usa_df.selectExpr("*").explain
1064
1065 DROP VIEW IF EXISTS just_usa_view;
1066
1067 69 database
1068 CREATE DATABASE some_db
1069 USE some_db
1070 SHOW tables
1071 SELECT current_database()
1072 DROP DATABASE IF EXISTS some_db;
1073
1074
1075 70 select
1076 SELECT
1077 CASE WHEN DEST_COUNTRY_NAME = 'UNITED STATES' THEN 1
1078 WHEN DEST_COUNTRY_NAME = 'Egypt' THEN 0
1079 ELSE -1 END
1080 FROM partitioned_flights
1081
1082 Global temp views are resolved regardless of database and are
1083 viewable across the entire Spark application, but they are removed at the end of the
1084 session
1085 CREATE GLOBAL TEMP VIEW just_usa_global_view_temp AS
1086 SELECT * FROM flights WHERE dest_country_name = 'United States'

```

```

1086
1087 71 complex type
1088 There are three core complex types in Spark SQL: structs, lists, and maps
1089
1090 Structs
1091 Structs are more akin to maps. They provide a way of creating or querying nested data
    in Spark.
1092 To create one, you simply need to wrap a set of columns (or expressions) in parentheses:
1093 CREATE VIEW IF NOT EXISTS nested_data AS
1094 SELECT (DEST_COUNTRY_NAME, ORIGIN_COUNTRY_NAME) as country, count FROM flights
1095 SELECT country.DEST_COUNTRY_NAME, count FROM nested_data
1096
1097
1098 List
1099 SELECT DEST_COUNTRY_NAME as new_name, collect_list(count) as flight_counts,
1100 collect_set(ORIGIN_COUNTRY_NAME) as origin_set
1101 FROM flights GROUP BY DEST_COUNTRY_NAME
1102
1103 sub query
1104 SELECT * FROM flights
1105 WHERE origin_country_name IN (SELECT dest_country_name FROM flights
1106 GROUP BY dest_country_name ORDER BY sum(count) DESC LIMIT 5)
1107
1108
1109 checking whether there is a flight that has the
1110 destination country as an origin and a flight that had the origin country as a
    destination:
1111 SELECT * FROM flights f1
1112 WHERE EXISTS (SELECT 1 FROM flights f2
1113 WHERE f1.dest_country_name = f2.origin_country_name)
1114 AND EXISTS (SELECT 1 FROM flights f2
1115 WHERE f2.dest_country_name = f1.origin_country_name)
1116
1117 spark.sql.shuffle.partitions 200 (default)
1118
1119
1120 71 cluster mode
1121 The cluster
1122 manager then launches the driver process on a worker node inside the cluster, in
    addition to the
1123 executor processes
1124
1125 cliet mode
1126 Spark driver remains on the client
1127 machine that submitted the application. This means that the client machine is
    responsible for
1128 maintaining the Spark driver process, and the cluster manager maintains the executor
    processes
1129
1130 Local mode
1131 it runs the entire Spark
1132 Application on a single machine
1133
1134 72 The Life Cycle of a Spark Application
1135 1) Client Request
1136 ./bin/spark-submit \
1137 --class <main-class> \
1138 --master <master-url> \
1139 --deploy-mode cluster \
1140 --conf <key>=<value> \
1141 ... # other options
1142 <application-jar> \
1143 [application-arguments]
1144
1145 2) Launch
1146 The SparkSession will subsequently communicate with the cluster manager

```

1147 (the darker line), asking it to launch Spark executor processes across the cluster  
1148 3)Execution  
1149 The driver and the workers communicate among themselves, executing code and  
1150 moving data around. The driver schedules tasks onto each worker, and each worker responds  
1151 with the status of those tasks and success or failure  
1152 4) After a Spark Application completes, the driver processs exits with either success  
or failure

1153  
1154 73  
1155 The Life Cycle of a Spark Application (Inside Spark)  
1156  
1157 1) session  
1158 val spark = SparkSession.builder().appName("Databricks Spark Example")  
1159 .config("spark.sql.warehouse.dir", "/user/hive/warehouse")  
1160 .getOrCreate()  
1161  
1162 SparkContext, you can create RDDs, accumulators, and broadcast variables  
1163 import org.apache.spark.SparkContext  
1164 val sc = SparkContext.getOrCreate()  
1165  
1166 It is important to note that you should never need to use the SQLContext and  
1167 rarely need to use the SparkContext  
1168  
1169 2) A Spark Job  
1170 In general, there should be one Spark job for one action. Actions always return  
results. Each job  
1171 breaks down into a series of stages, the number of which depends on how many shuffle  
1172 operations need to take place  
1173  
1174 stages  
1175 In general, Spark will try to pack as much work as possible (i.e.,  
1176 as many transformations as possible inside your job) into the same stage, but the  
engine starts  
1177 new stages after operations called shuffles  
1178 c  
1179 A shuffle represents a physical repartitioning of the  
1180 data—for example, sorting a DataFrame, or grouping data that was loaded from a file  
1181  
1182 Spark starts a new stage after each  
1183 shuffle  
1184  
1185 The  
1186 spark.sql.shuffle.partitions default value is 200, which means that when there is a  
shuffle  
1187 performed during execution, it outputs 200 shuffle partitions by default  
1188  
1189 entire stage is computed in parallel. The final result aggregates those  
1190 partitions individually, brings them all to a single partition before finally sending  
the final result  
1191 to the driver  
1192  
1193  
1194 Task:  
1195 A task is just a unit of computation applied to a unit of data (the  
1196 partition)  
1197  
1198 74 shuffle  
1199 Spark always executes  
1200 shuffles by first having the "source" tasks (those sending data) write shuffle files to  
their local  
1201 disks during their execution stage. Then, the stage that does the grouping and  
reduction launches  
1202 and runs tasks that fetch their corresponding records from each shuffle file and  
performs that  
1203 computation  
1204

```
1205 75 spark-submit
1206 --packages Comma-separated list of Maven coordinates of JARs to include on the driver
      and executor
1207 classpaths
1208 --drivermemory
1209 --executormemory
1210 MEM
1211 --driverjava-
1212 options
1213
1214 ./bin/spark-submit \
1215 --class org.apache.spark.examples.SparkPi \
1216 --master spark://207.184.161.138:7077 \
1217 --executor-memory 20G \
1218 --total-executor-cores 100 \
1219 replace/with/path/to/examples.jar \
1220 1000
1221
1222
1223 76
1224 The SparkConf
1225
1226 val conf = new SparkConf().setMaster("local[2]").setAppName("DefinitiveGuide")
1227 .set("some.conf", "to.some.value")
1228
1229 ./bin/spark-submit --name "DefinitiveGuide" --master local[4] ...
1230
1231 77.
1232 Application Properties
1233 spark.driver.cores 1 Number of cores to use for the driver process, only in cluster mode.
1234 spark.driver.memory 1g
1235 spark.executor.memory 1g Amount of memory to use per executor process (e.g., 2g, 8g).
1236 spark.submit.deployMode (none) client/cluster
1237
1238 spark.files.maxPartitionBytes (maximum
1239 partition size when reading files)
1240
1241 78Configuring Memory Management
1242
1243 79 scheduler
1244 By default, Spark's scheduler runs jobs in FIFO fashion
1245 Under fair sharing, Spark assigns tasks
1246 between jobs in a round-robin fashion so that all jobs get a roughly equal share of
      cluster
1247 resources spark.scheduler.mode property to FAIR
1248
1249 80 cluster
1250 onpremises
1251 deployment, your cluster is fixed in size
1252 on-premises, the best way to combat the resource utilization problem
1253 is to use a cluster manager that allows you to run many Spark applications and
      dynamically
1254 reassign resources between them
1255
1256 a better idea to use global storage systems that are decoupled
1257 from a specific cluster, such as Amazon S3, Azure Blob Storage, or Google Cloud Storage
      and
1258 spin up machines dynamically for each Spark workload
1259
1260
1261 Standalone Mode
1262 disadvantage of the standalone mode is that it's more limited than the other cluster
      managers—in
1263 particular, your cluster can only run Spark
1264
1265 The first step is to start the master process on the machine
```

1266 \$SPARK\_HOME/sbin/start-master.sh  
1267 Once started, the master prints out a `spark://HOST:PORT` URI  
1268  
1269 With that URI, start the worker nodes by  
1270 logging in to each machine and running the following script using the URI you just  
1271 received  
1272 from the master node  
1273 \$SPARK\_HOME/sbin/start-slave.sh <master-spark-URI>  
1274  
81  
1275 Hadoop configurations  
1276 If you plan to read and write from HDFS using Spark, you need to include two Hadoop  
1277 configuration files on Spark's classpath: `hdfs-site.xml`, which provides default  
1278 behaviors for the  
1279 HDFS client; and `core-site.xml`  
1280  
1281 To make these files visible to Spark, set `HADOOP_CONF_DIR` in `$SPARK_HOME/spark-env.sh`  
1282 to a  
1283 location containing the configuration files or as an environment variable  
1284  
82  
1285 Mesos is the heaviest-weight cluster manager, simply because you might  
1286 choose this cluster manager only if your organization already has a large-scale  
1287 deployment of  
1288 Mesos  
1289  
1289 Coarse-grained mode means that each Spark executor runs as a single  
1290 Mesos task  
1291  
1292 `val spark = SparkSession.builder`  
1293 `.master("mesos://HOST:5050")`  
1294 `.appName("my app")`  
1295 `.config("spark.executor.uri", "<path to spark-2.2.0.tar.gz uploaded above>")`  
1296 `.getOrCreate()`  
1297  
1298 Dynamic allocation  
1299 dynamic allocation (described next) can be turned on to let applications  
1300 scale up and down dynamically based on their current number of pending tasks  
1301 set `spark.dynamicAllocation.enabled` to `true`,  
1302 default is disabled.  
1303  
1304 second:  
1305 `spark.shuffle.service.enabled` to `true` in your application. The purpose of  
1306 the external shuffle service is to allow executors to be removed without deleting  
1307 shuffle files  
1308 written by them  
1309  
83,  
1310 YARN is great for HDFS-based applications but is not commonly used for  
1311 much else. Additionally, it's not well designed to support the cloud  
1312  
1313 Spark standalone mode is the lightest-weight cluster  
1314 manager and is relatively simple to understand and take advantage of, but then you're  
1315 going to  
1316 be building more application management infrastructure  
1317  
1317 it doesn't really make sense to have a Mesos cluster  
1318 for only running Spark Applications  
1319  
84  
1320 The metrics system is  
1321 configured via a configuration file that Spark expects to be present at  
1322 `$SPARK_HOME/conf/metrics.properties`  
1323 These metrics can be output to a variety of  
1324 different sinks, including cluster monitoring solutions like Ganglia

1326  
1327 spark.sparkContext.setLogLevel("INFO")  
1328 The logs themselves will be printed to standard error when running a local  
1329 mode application, or saved to files by your cluster manager when running Spark on a  
cluster.  
1330 Refer to each cluster manager's documentation about how to find them—typically, they are  
1331 available through the cluster manager's web UI  
1332  
1333 86 spark GUI  
1334  
1335 If  
1336 you're running multiple applications, they will launch web UIs on increasing port numbers  
1337 (4041, 4042, ...)  
1338  
1339  
1340 87 Error troubleshooting  
1341  
1342 Ensure that machines can communicate with one another  
1343 Ensure that your Spark resource configurations are correct and that your cluster manager  
1344 is properly set up for Spark  
1345 Double-check to verify that the cluster has the network connectivity  
1346 There might be issues with libraries or classpaths that are causing the wrong version  
of a  
1347 library to be loaded for accessing storage  
1348  
1349 88  
1350 Potential treatments  
1351 Check to see if your data exists or is in the format that you expect  
1352  
1353 If an error quickly pops up when you run a query (i.e., before tasks are launched), it is  
1354 most likely an analysis error while planning the query  
1355  
1356 It's also possible that your own code for processing the data is crashing, in which case  
1357 Spark will show you the exception thrown by your code. In this case, you will see a task  
1358 marked as "failed" on the Spark UI,  
1359  
1360  
1361 89 Slow Tasks or Stragglers  
1362 symptoms:  
1363 Spark stages seem to execute until there are only a handful of tasks left. Those tasks  
1364 then take a long time.  
1365 These slow tasks show up in the Spark UI and occur consistently on the same dataset(s).  
1366 These occur in stages, one after the other.  
1367 Scaling up the number of machines given to the Spark Application doesn't really help—  
1368 some tasks still take much longer than others.  
1369 In the Spark metrics, certain executors are reading and writing much more data than  
1370 others  
1371  
1372 solution:  
1373 Try increasing the number of partitions to have less data per partition  
1374 Try repartitioning by another combination of columns  
1375 Check whether your user-defined functions (UDFs) are wasteful in their object  
1376 allocation or business logic. Try to convert them to DataFrame code if possible  
1377 Try increasing the memory allocated to your executors if possible  
1378 Monitor the executor that is having trouble  
1379 Another common issue can arise when you're working with Datasets. Because Datasets  
1380 perform a lot of object instantiation to convert records to Java objects for UDFs, they  
1381 can cause a lot of garbage collection  
1382  
1383 90 Slow Aggregations  
1384 Slow tasks during a groupBy call. Jobs after the aggregation are slow, as well  
1385  
1386 treatments:  
1387 Increasing the number of partitions, prior to an aggregation  
1388 Increasing executor memory can help alleviate this issue, as well. If a single key has  
lots

1389 of data, this will allow its executor to spill to disk less often and finish faster,  
1390  
1391 If you find that tasks after the aggregation are also slow, this means that your dataset  
1392 might have remained unbalanced after the aggregation. Try inserting a repartition  
1393 call to partition it randomly  
1394  
1395 Ensuring that all filters and SELECT statements that can be are above the aggregation can  
1396 help to ensure that you're working only on the data that you need to be working  
1397 Ensure null values are represented correctly (using Spark's concept of null) and not as  
1398 some default value like " " or "EMPTY". Spark often optimizes for skipping nulls  
1399  
1400 For instance,  
1401 collect\_list and collect\_set are very slow aggregation functions because they  
1402 must return all the matching objects to the driver  
1403  
1404 91  
1405 Slow Joins  
1406 A join stage seems to be taking a long time , Stages before and after the join seem to  
be operating normally  
1407  
1408 treatments:  
1409 Many joins can be optimized  
1410  
1411 Experimenting with different join orderings can really help speed up jobs, especially if  
1412 some of those joins filter out a large amount of data; do those first.  
1413  
1414 Partitioning a dataset prior to joining can be very helpful for reducing data movement  
1415 across the cluster, especially if the same dataset will be used in multiple join  
operations.  
1416 It's worth experimenting with different prejoin partitioning. Keep in mind, again, that  
1417 this isn't "free" and does come at the cost of a shuffle.  
1418  
1419 Ensuring that all filters and select statements that can be are above the join  
1420  
1421 Ensure that null values are handled correctly  
1422  
1423 Sometimes Spark can't properly plan for a broadcast join if it doesn't know any  
1424 statistics about the input DataFrame or table. If you know that one of the tables that  
you  
1425 are joining is small, you can try to force a broadcast  
1426  
1427 92 Slow Reads and Writes  
1428 Slow reading of data from a distributed file system  
1429  
1430 treatments:  
1431 Turning on speculation (set spark.speculation to true) can help with slow reads and  
1432 writes. This will launch additional tasks with the same operation in an attempt to see  
1433 whether it's just some transient issue in the first task  
1434  
1435 Ensuring sufficient network connectivity can be important  
1436 For distributed file systems such as HDFS running on the same nodes as Spark, make  
1437 sure Spark sees the same hostnames for nodes as the file system  
1438  
1439 93 Driver OutOfMemoryError or Driver Unresponsive  
1440 Signs and symptoms  
1441 Spark Application is unresponsive or crashed.  
1442 OutOfMemoryErrors or garbage collection messages in the driver logs.  
1443 Commands take a very long time to run or don't run at all.  
1444 Interactivity is very low or non-existent.  
1445 Memory usage is high for the driver JVM.  
1446  
1447 Potential treatments  
1448 Your code might have tried to collect an overly large dataset to the driver node using  
1449 operations such as collect  
1450 You might be using a broadcast join where the data to be broadcast is too big  
1451

1452 A long-running application generated a large number of objects on the driver and is  
1453 unable to release them. Java's jmap tool can be useful to see what objects are filling  
1454 most of the memory  
1455  
1456 Issues with JVMs running out of memory can happen if you are using another language  
1457 binding, such as Python, due to data conversion between the two requiring too much  
1458 memory in the JVM  
1459  
1460 94 Executor OutOfMemoryError or Executor Unresponsive  
1461 OutOfMemoryErrors or garbage collection messages in the executor logs. You can find  
1462 these in the Spark UI.  
1463 Executors that crash or become unresponsive.  
1464 Slow tasks on certain nodes that never seem to recover  
1465  
1466 Try increasing the memory available to executors and the number of executors.  
1467 Try increasing PySpark worker size via the relevant Python configurations.  
1468 Look for garbage collection error messages in the executor logs. Some of the tasks that  
1469 are running, especially if you're using UDFs, can be creating lots of objects that need  
to  
1470 be garbage collected.  
1471  
1472 95  
1473 No Space Left on Disk Errors  
1474 If you have a cluster with limited storage space, some nodes may run out first due to  
1475 skew. Repartitioning the data as described earlier may help here.  
1476  
1477 Some of these determine how long logs should be kept on the machine  
1478  
1479 Serialization Errors  
1480 Try not to refer to any fields of the enclosing object in your UDFs when creating UDFs  
1481 inside a Java or Scala class. This can cause Spark to try to serialize the whole  
enclosing  
1482 object, which may not be possible  
1483  
1484  
1485  
1486  
1487