

Текстовая расшифровка видео:

СМЕЩЕННЫЕ ДАННЫЕ

План:

- Равномерно ли распределены данные;
- Data Skew;
- Виды партиционирования – RoundRobin;
- Виды партиционирования – HashPartitioning;
- Другой вариант – coalesce;
- «Соленые» данные;
- Подытог.

Равномерно ли распределены данные?

Одновременно Spark выполняет N task'ов, которые обрабатывают N партиций, где N – это суммарное число доступных потоков на всех воркерах.

Исходя из этого, важно обеспечивать:

- Достаточное количество партиций для распределения нагрузки по всем воркерам;
- Равномерное распределение данных между партициями.

Data Skew

Предположим, что есть датасет, который крайне плохо распределен по кластеру:



▶ Запустить стенд



Дедлайн 07 июля, 23:59 Мск



```

1. from pyspark.sql.functions import when, col, lit
2.
3. def getItemsPerPartition(ds):
4.     def count_in_a_partition(idx, iterator):
5.         yield idx, sum(1 for _ in iterator)
6.     return ds.rdd.mapPartitionsWithIndex(count_in_a_partition).collect()
7.
8.
9. skew_column = when(col("id") < 900, lit(0)).otherwise(lit(1))
10. skewDF = spark.range(0,1000).repartition(10, skew_column)
11. getItemsPerPartition(skewDF)

```

У нас есть и Spark, и функции. Мы делаем dataframe, вызываем repartition и говорим «сделай 10 партиций, но при распределении используй эту колонку»:

```

[1]: from pyspark import SparkContext, SparkConf
    from pyspark.sql import SparkSession

    conf = SparkConf().setAppName("my-app").setMaster("local[4]")
    sc = SparkContext(conf=conf)
    spark = SparkSession(sc)

[2]: from pyspark.sql.functions import when, col, lit, spark_partition_id, collect_list

[3]: skew_column = when(col("id") < 900, lit(0)).otherwise(lit(1))

[ ]: skewDF = spark.range(0,1000).repartition(10, skew_column)

[ ]: def getItemsPerPartition(ds):
    def count_in_a_partition(idx, iterator):
        yield idx, sum(1 for _ in iterator)
    return ds.rdd.mapPartitionsWithIndex(count_in_a_partition).collect()

[ ]: getItemsPerPartition(skewDF)

[ ]: # здесь мы передаем только новое количество партиций и Spark выполнит RoundRobinPartitioning
    rrDF = skewDF.repartition(10)
    getItemsPerPartition(rrDF)

```

Часто там может быть ID и т.д.

Если посмотреть на то, по какому условию мы делаем партиционирование, то увидим, что для чисел от 0 до 1000 значения меньше 900 попадут в одну партицию, а остальные – в другую, так как в колонке два значения. Чувствуется, что уже что-то не так.

Есть функция, которая называется «**mapPartitionsWithIndex**». Мы проходимся по всем партициям и смотрим, сколько записей в каждой партиции оказалось. Мы видим, что несмотря на то, что партиций 10, данные оказались в двух:

```

[2]: from pyspark.sql.functions import when, col, lit, spark_partition_id, collect_list

[3]: skew_column = when(col("id") < 900, lit(0)).otherwise(lit(1))

[4]: skewDF = spark.range(0,1000).repartition(10, skew_column)

[5]: def getItemsPerPartition(ds):
    def count_in_a_partition(idx, iterator):
        yield idx, sum(1 for _ in iterator)
    return ds.rdd.mapPartitionsWithIndex(count_in_a_partition).collect()

[6]: getItemsPerPartition(skewDF)

[6]: [(0, 0),
      (1, 900),
      (2, 0),
      (3, 100),
      (4, 0),
      (5, 0),
      (6, 0),
      (7, 0),
      (8, 0),
      (9, 0)]

[ ]: # здесь мы передаем только новое количество партиций и Spark выполнит RoundRobinPartitioning

```

Причем в партиции «1» окажется 900 записей, а в партиции «3» – 100. Это крайне нежелательный случай!

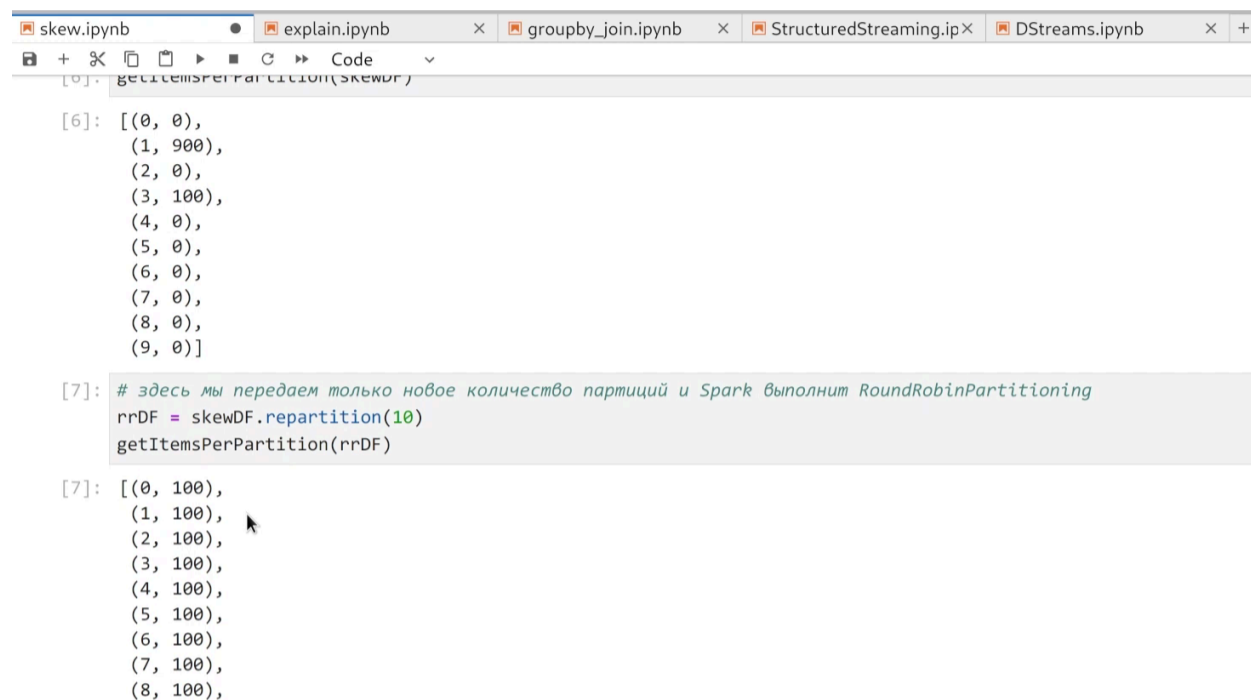
Виды партиционирования – RoundRobin

```
# здесь мы передаем только новое
# количество партиций и Spark выполнит
# RoundRobinPartitioning
rrDF = skewDF.repartition(10)
getItemsPerPartition(rrDF)
```

RoundRobinPartitioning – это подход, при котором Spark просматривает все значения в партициях и каждое следующее значение перемещает на следующую партицию.

RoundRobinPartitioning может быть эффективен, но **он ничего не знает о данных**. Он пересылает **очень** много данных горизонтально.

Выполним код:



```
skew.ipynb explain.ipynb groupby_join.ipynb StructuredStreaming.ipynb DStreams.ipynb
Code
[6]: [(0, 0),
      (1, 900),
      (2, 0),
      (3, 100),
      (4, 0),
      (5, 0),
      (6, 0),
      (7, 0),
      (8, 0),
      (9, 0)]

[7]: # здесь мы передаем только новое количество партиций и Spark выполнит RoundRobinPartitioning
      rrDF = skewDF.repartition(10)
      getItemsPerPartition(rrDF)

[7]: [(0, 100),
      (1, 100),
      (2, 100),
      (3, 100),
      (4, 100),
      (5, 100),
      (6, 100),
      (7, 100),
      (8, 100),
```

RoundRobin'ом данные раскидались.

Виды партиционирования – HashPartitioning

Как мы говорили ранее, имя колонки определяется по значениям. На самом деле там используется тип Partitioning'a, который называется «**HashPartitioning**». От значений, которые есть в колонке, считается хэш. На основе этого хэша принимается решение в какую партицию какие данные закинуть:

```
# здесь мы добавляем к числу партиций колонку, по которой необходимо
# сделать репартиционирование, поэтому Spark выполнит
# HashPartitioning
hpDF = skewDF.repartition(10, col("id"))
getItemsPerPartition(hpDF)
```

```
[(0, 93),
 (1, 107),
 (2, 93),
 (3, 105),
 (4, 96),
 (5, 118),
 (6, 89),
 (7, 97),
 (8, 95),
 (9, 107)]
```

Если мы запустим партиционирование по этой колонке, указав ее, то выполнится HashPartitioning.

Иногда могут быть коллизии, и распределение данных по кластеру будет не таким равномерным. Если мы знаем, что кусочки датасетов, относятся к одному и тому же значению в этой колонке, будем процессить независимо. Подобный метод партиционирования может быть более эффективным при расчетах.

```

skew.ipynb explain.ipynb groupby_join.ipynb StructuredStreaming.ipynb DStreams.ipynb
Code
(8, 100),
(9, 100)]

[8]: # здесь мы добавляем к числу партиций колонку, по которой необходимо
# сделать репартиционирование, поэтому Spark выполнит
# HashPartitioning
hpDF = skewDF.repartition(10, col("id"))
getItemsPerPartition(hpDF)

[8]: [(0, 93),
(1, 107),
(2, 93),
(3, 105),
(4, 96),
(5, 118),
(6, 89),
(7, 97),
(8, 95),
(9, 107)]

[ ]: df1 = spark.range(0, 25).repartition(6)
df1.groupBy(spark_partition_id()).agg(collect_list("id")).show()

df2 = df1.repartition(4)
df2.groupBy(spark_partition_id()).agg(collect_list("id")).show()

```

Другой вариант – coalesce

Когда данные распартиционированы, мы хотим их переложить. Это можно сделать двумя способами:

- **Repartition**, когда речь идет об увеличении партиций.
- **Coalesce**, когда речь идет об уменьшении партиций.

```

df1 = spark.range(0, 25).repartition(6)
df1.groupBy(spark_partition_id()).agg(collect_list("id")).show()

df2 = df1.repartition(4)
df2.groupBy(spark_partition_id()).agg(collect_list("id")).show()

df3 = df1.coalesce(4)
df3.groupBy(spark_partition_id()).agg(collect_list("id")).show()

```

Coalesce пытается минимизировать дополнительные пересылки данных на кластере.

Рассмотрим пример:

```

skew.ipynb explain.ipynb groupby_join.ipynb StructuredStreaming.ipynb DStreams.ipynb
Code
(9, 107)]

[9]: df1 = spark.range(0, 25).repartition(6)
df1.groupBy(spark_partition_id()).agg(collect_list("id")).show()

df2 = df1.repartition(4)
df2.groupBy(spark_partition_id()).agg(collect_list("id")).show()

df3 = df1.coalesce(4)
df3.groupBy(spark_partition_id()).agg(collect_list("id")).show()

+-----+-----+
|SPARK_PARTITION_ID()| collect_list(id)|
+-----+-----+
|                0|[2, 10, 13, 19, 24]|
|                1|[4, 6, 15, 22]|
|                2|[1, 8, 16, 20]|
|                3|[3, 9, 12, 23]|
|                4|[0, 11, 17, 21]|
|                5|[5, 7, 14, 18]|
+-----+-----+

+-----+-----+
|SPARK_PARTITION_ID()| collect_list(id)|
+-----+-----+

```

Мы видим данные, раскиданные по шести партициям. Вызываем базовый Repartition. Как мы помним, происходит RoundRobinPartitioning, и данные перекадываются по всему кластеру как получится. Если же вызвать Coalesce, мы увидим, что партиций стало меньше, но при этом структура партиций осталась. Некоторые сегменты просто объединились между собой. Например, две партиции лежат на одном воркере, а мы уменьшаем количество партиций. При вызове Coalesce есть серьезный шанс схлопнуть две партиции в одну, при этом дополнительной пересылки по кластеру не будет.

Нюанс: после выполнения Coalesce количество данных по партициям будет неравномерным.

«Соленые» данные

```
airports = spark.read.options(header=True, inferSchema=True).csv(
    os.path.expanduser("~/Experiments/slurm/airport-codes.csv")
)
airports.groupBy(col("type")).count().orderBy(col("count").desc()).show()
```

Мы должны получить следующую таблицу:

```
+-----+-----+
|      type|count|
+-----+-----+
| small_airport|34808|
|      heliport|12028|
| medium_airport| 4537|
|      closed| 4378|
| seaplane_base| 1030|
| large_airport|  616|
| balloonport|   24|
+-----+-----+
```

Все получилось:

```
[12]: airports = spark.read.options(header=True, inferSchema=True).csv(
    os.path.expanduser("~/Experiments/slurm/airport-codes.csv")
)
airports.groupBy(col("type")).count().orderBy(col("count").desc()).show()

+-----+-----+
|      type|count|
+-----+-----+
| small_airport|34808|
|      heliport|12028|
| medium_airport| 4537|
|      closed| 4378|
| seaplane_base| 1030|
| large_airport|  616|
| balloonport|   24|
+-----+-----+
```

В ней мы видим **нюанс**: мы проводим группировку по колонке с типом. В колонке с типом большой разброс по количеству уникальных значений, например, маленьких аэропортов $\approx 35\ 000$, а аэропортов для воздушных шаров – 24. Когда происходит подобное, есть риск, что при выполнении GroupBy, Spark может вылететь по памяти. Внутри происходит shuffle, и Spark репроцессионирует данные, пересылая данные, относящиеся к одному значению колонки, на одну ноду. Это одна из самых частых ошибок.

Для решения такой проблемы используется **подход «солоние данных»**. Попробуем посчитать ту же задачу. Мы сделаем не одну группировку, а две.

Для начала придумаем колонку, в которой зададим распределение. Здесь примерно прикидываем:

```
[13]: from pyspark.sql.functions import pmod, round, rand, sum

[ ]: # дополнительная колонка с солью
salt_mod_ten = pmod(round((rand() * 100), 0), lit(10)).cast("int")
salted = airports.withColumn("salt", salt_mod_ten)
# агрегируем более взвешенно
first_step = salted.groupBy(col("type"), col("salt")).count()
second_step = first_step.groupBy(col("type")).agg(sum("count").alias("count"))
second_step.show()

[ ]:
```

Например, у нас будет 10 партиций. Мы заполняем колонку остатками от деления. Рандомное число умножаем на 100, округляем его до целого и получаем какое-то значение. Считаем остаток от деления на 10. У нас получается выдуманная «солонная» колонка.

Смотрим:

```
|ident|      type|      name|elevation_ft|continent|iso_country|iso_region|municipality|gps_code|iata_code|local_code|
coordinates|salt|
+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
| 00A| heliport| Total Rf Heliport|      11|      NA|      US|      US-PA| Bensalem| 00A| null| 00A|-
74.9336013793945...| 8|
| 00AA|small_airport|Aero B Ranch Airport|     3435|      NA|      US|      US-KS| Leoti| 00AA| null| 00AA|-
101.473911, 38.7...| 3|
| 00AK|small_airport| Lowell Field|      450|      NA|      US|      US-AK|Anchor Point| 00AK| null| 00AK|-
151.695999146, 5...| 4|
| 00AL|small_airport| Epps Airpark|      820|      NA|      US|      US-AL| Harvest| 00AL| null| 00AL|-
86.7703018188476...| 0|
| 00AR| closed|Newport Hospital ...|     237|      NA|      US|      US-AR| Newport| null| null| null|-
-91.254898, 35.6087| 9|
| 00AS|small_airport| Fulton Airport|     1100|      NA|      US|      US-OK| Alex| 00AS| null| 00AS|-
97.8180194, 34.9...| 2|
| 00AZ|small_airport| Cordes Airport|     3810|      NA|      US|      US-AZ| Cordes| 00AZ| null| 00AZ|-
112.165000915527...| 1|
| 00CA|small_airport|Goldstone /Gts/ A...|     3038|      NA|      US|      US-CA| Barstow| 00CA| null| 00CA|-
116.888000488, 3...| 4|
| 00CL|small_airport|Williams A...|      87|      NA|      US|      US-CA| Bisco| 00CL| null| 00CL|-
116.888000488, 3...| 4|
```

В этой колонке будут числа от 0 до 10. Сделаем первый groupBy не по колонке «Type», а по двум колонкам: «Type», «Salt». Это приведет к тому, что данные по кластеру будут группироваться не

```
[13]: from pyspark.sql.functions import pmod, round, rand, sum
```

```
•[14... # дополнительная колонка с солью
salt_mod_ten = pmod(round((rand() * 100), 0), lit(10)).cast("int")
salted = airports.withColumn("salt", salt_mod_ten)
# агрегируем более взвешенно
first_step = salted.groupBy(col("type"), col("salt")).count()
#second_step = first_step.groupBy(col("type")).agg(sum("count").alias("count"))
#second_step.show()
```

Вторым шагом мы сделаем второй groupBy, в котором уже сможем применить агрегат к колонке «Type». Поскольку данные связаны с конкретным типом, они уже независимо лежат на партиции. Есть серьезный шанс, что это увидит оптимизатор, вследствие чего дополнительной второй пересылки не будет. Мы вызываем ту же операцию, которую считали выше:

```
# дополнительная колонка с солью
salt_mod_ten = pmod(round((rand() * 100), 0), lit(10)).cast("int")
salted = airports.withColumn("salt", salt_mod_ten)
# агрегируем более взвешенно
first_step = salted.groupBy(col("type"), col("salt")).count()
second_step = first_step.groupBy(col("type")).agg(sum("count").alias("count"))
#second_step.show()
```

Получаем тот же результат, но из-за группировки в два этапа мы не вылетим по памяти:

```
[15]: # дополнительная колонка с солью
salt_mod_ten = pmod(round((rand() * 100), 0), lit(10)).cast("int")
salted = airports.withColumn("salt", salt_mod_ten)
# агрегируем более взвешенно
first_step = salted.groupBy(col("type"), col("salt")).count()
second_step = first_step.groupBy(col("type")).agg(sum("count").alias("count"))
second_step.show()
```

```
+-----+
|      type|count|
+-----+
| large_airport| 616|
|  balloonport|  24|
| seaplane_base|1030|
|   heliport|12028|
|   closed| 4378|
|medium_airport| 4537|
| small_airport|34808|
+-----+
```

Подытог

Партиционирование – важный аспект распределенных вычислений, от которого напрямую зависит стабильность и скорость вычислений.

В Spark работает правило **1 TASK=1 THREAD= 1 PARTITION**

Репартиционирование и соление данных позволяет решить проблему перекоса данных и вычислений.

Важно помнить, что репартиционирование использует дисковую и сетевую подсистемы – обмен данными происходит по сети, а результат часто записывается на диск, что может стать узким местом при выполнении репартиционирования.

Как вам урок?



Изучил, далее >

Слёрм ©

+7 (495) 248-05-80

[Лицензия №ДЛ-1368 от 22.08.2019](#)

[Политика конфиденциальности](#)

[Публичная оферта](#)