

spark access first n rows - take vs limit

Asked 7 years ago Modified 1 year, 1 month ago Viewed 169k times



I want to access the first 100 rows of a spark data frame and write the result back to a CSV file.

54

Why is `take(100)` basically instant, whereas



```
df.limit(100)
  .repartition(1)
  .write
  .mode(SaveMode.Overwrite)
  .option("header", true)
  .option("delimiter", ";")
  .csv("myPath")
```



takes forever. I do not want to obtain the first 100 records per partition but just any 100 records.

Why is `take()` so much faster than `limit()`?

apache-spark

apache-spark-sql

limit

Share Edit Follow

edited Feb 7, 2021 at 17:11



Marioanzas

1,895 ● 2 ● 18 ● 36

asked Oct 19, 2017 at 14:31



Georg Heiler

17.6k ● 40 ● 172 ● 311

1 So, you can use `take(100)` indeed; what is the question? – [desertnaut](#) Oct 20, 2017 at 10:39 ✎

2 Why take is so much faster than limit. – [Georg Heiler](#) Oct 20, 2017 at 11:10 ✎

1 Indeed I could, but so far have not seen a way to create a df of the local array to use Sparks nice CSV handling capabilities. Limit should just provide this. – [Georg Heiler](#) Oct 20, 2017 at 11:13

6 Answers

Sorted by: Highest score (default) ▾



Although it still is answered, I want to share what I learned.

68

```
myDataFrame.take(10)
```



-> **results in an Array of Rows.** This is an action and performs collecting the data (like collect does).





```
myDataFrame.limit(10)
```



-> **results in a new Dataframe.** This is a transformation and does not perform collecting the data.

I do not have an explanation why then limit takes longer, but this may have been answered above. This is just a basic answer to what the difference is between take and limit.

Share Edit Follow

edited Oct 15, 2020 at 9:23

answered Jun 26, 2019 at 12:09



pfnuesel

15.2k ● 14 ● 60 ● 73



Kaspato

1,353 ● 2 ● 13 ● 28

The difference between action and transformation is correct, but that does not explain why limit should take longer than take (once the plan executes). – [Arjen P. De Vries](#) Nov 11, 2020 at 8:16



26



This is because predicate pushdown is currently not supported in Spark, see [this very good answer](#).

Actually, take(n) should take a really long time as well. I just tested it, however, and get the same results as you do - take is almost instantaneous irregardless of database size, while limit takes a lot of time.

Share Edit Follow

edited Mar 16, 2018 at 9:45

answered Mar 16, 2018 at 9:14



Thomas

5,064 ● 5 ● 44 ● 74

Collect only works in spark dataframes. When I collect first 100 rows it is instant and data resides in memory as a regular list. Collect in sparks sense is then no longer possible. – [Georg Heiler](#) Mar 16, 2018 at 9:35

You are right of course, I forgot take returns a list. I just tested it, and get the same results - I expected both take and limit to be slow. – [Thomas](#) Mar 16, 2018 at 9:47

stackoverflow.com/questions/35869884/... <- This question however explicitly states that others have problems with `take()` as well - which version of pyspark are you using? – [Thomas](#) Mar 16, 2018 at 9:48

1 Spark scala 2.2 – [Georg Heiler](#) Mar 16, 2018 at 11:32



limit work for me:

-1

```
limitDF= df.limit(5)
```



Even better approach is limit your data with filter



```
filterDF= df.filter("name = 'jitendra']").limit(5)
```



Share Edit Follow

answered Oct 4, 2023 at 5:43

**Jitendra Bansiwal**

113 ● 1 ● 4



Limit() will not work in partition, so it will take more time to execute

-1

Share Edit Follow

answered Feb 18, 2022 at 14:12

**Ratheesh**

19 ● 5



You can use take(n) to limit the data. Adding the complete code with output in the screenshot.

-3

```

1 import org.apache.spark.{SparkConf, SparkContext}
2 import org.apache.spark.sql.{SparkSession}
3
4 object SparkSessionTest extends App{
5   val spark = SparkSession.builder()
6     .master("local[3]")
7     .appName("SparkByExample")
8     .getOrCreate()
9   spark.sparkContext.setLogLevel("ERROR")
10
11   val rdd3=spark.sparkContext.textFile("C:/tmp/abc.csv")
12   val rdd6 = rdd3.map(f=>f.split(",")).take(10)
13   rdd6.foreach(f=>println(f(0), f(1)))
14   println(rdd6)
15 }

```

```

21/03/15 12:35:37 INFO SharedState: Setting hive.metastore.warehouse.dir ('null') to the value of spark.sql.warehouse.dir ('file:/C:/Users/shiv/IdeaProjects/SparkHadoop/spark-warehouse/').
21/03/15 12:35:37 INFO SharedState: Warehouse path is 'file:/C:/Users/shiv/IdeaProjects/SparkHadoop/spark-warehouse/'.
21/03/15 12:35:39 INFO StateStoreCoordinatorRef: Registered StateStoreCoordinator endpoint
(Year,Industry_aggregation_NZSIOC)
(2019,Level 1)
(2019,Level 1)
(2019,Level 1)
(2019,Level 1)
(2019,Level 1)
(2019,Level 1)
(2019,Level 1)
(2019,Level 1)
(2019,Level 1)
(2019,Level 1)
[[Ljava.lang.String;@5a49af50
Process finished with exit code 0

```

Share Edit Follow

answered Mar 15, 2021 at 7:08

**Shyam Gupta**

503 ● 4 ● 8



.take() could be the answer, but I used a simple head command like below

-9

df.head(3)

.take() did not work for me.

[Share](#) [Edit](#) [Follow](#)[edited Jan 7, 2021 at 17:59](#)[answered Jan 7, 2021 at 12:35](#)**10 Rep****2,270** ● 7 ● 20 ● 33**Vignesh M21****1**

1 Unfortunately, this is not an answer, it's more a feeling ... – [Gilles Bodart](#) Oct 18, 2022 at 12:59
