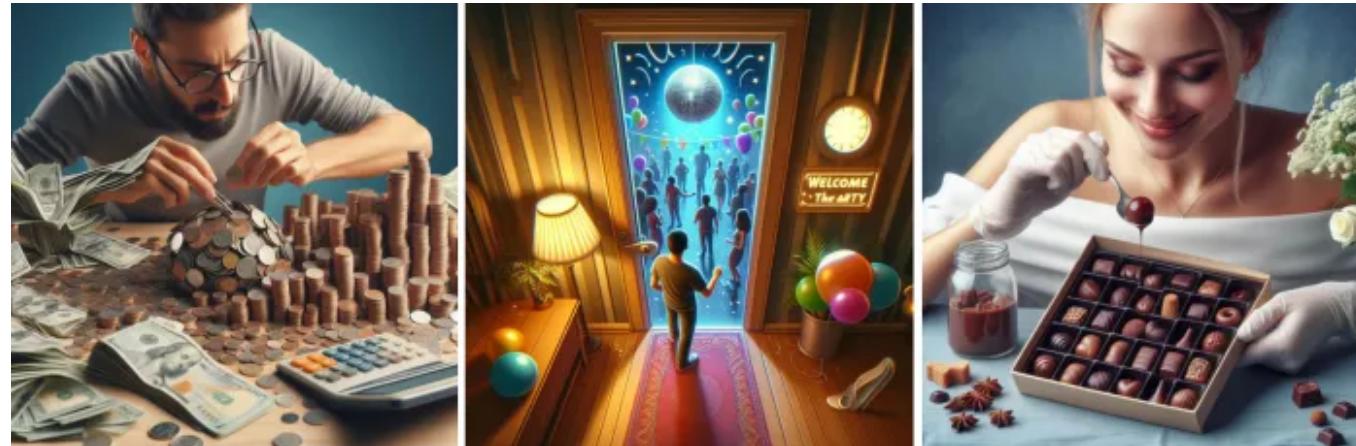


Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#)`count (), isEmpty (), and take (n)`

Optimizing Apache Spark: Strategic Use of `count()`, `isEmpty()`, and `take(n)`

Asish Panda · [Follow](#)

7 min read · Mar 9, 2024



Introduction

Discover the secrets of efficient data processing with Apache Spark in our latest exploration. This blog post dives deep into the practical applications and performance considerations of three essential Apache Spark actions: `count()`, `isEmpty()`, and `take(n)`. As the go-to engine for big data analytics, Apache Spark's capabilities are unparalleled. Yet, to leverage its full potential, developers and data engineers must master these actions to optimize their data processing workflows effectively. Join us as we unravel the intricacies of these functions, offering you a roadmap to enhanced performance in handling vast datasets. In Apache Spark, an **action** is a type of operation that triggers the execution of the computation instructions specified in your Spark program. Unlike transformations, which are lazy and only define a new RDD or DataFrame without triggering any computation, actions force the computation to be evaluated and results to be produced. In Apache Spark, `count()`, `isEmpty()`, and `take(n)` are one of the few different action methods used for different purposes when working with DataFrames. Understanding the context and performance implications of each can help you decide which is best to use and when.

Audience for this Blog

Targeted at data engineers, Spark developers, and data scientists, this blog demystifies the complexities of Apache Spark for those entrenched in the world of large-scale data. Whether you're looking to enhance the efficiency of your Spark jobs or make strategic decisions based on dataset characteristics, this guide serves as your indispensable resource for mastering the nuances of Spark's most crucial actions.

Prerequisites

A foundational understanding of Apache Spark lays the groundwork for this guide. Readers should be familiar with Spark's architecture and foundational concepts, such as RDDs (Resilient Distributed Datasets) and DataFrames. Additionally, a grasp of Spark's transformation and action operations will enrich your learning experience, setting the stage for advanced data manipulation and analysis techniques.

Problem Statement

In the realm of big data, efficiency and speed are kings. Processing data at scale requires not just understanding the dataset's size and content but also adopting the most effective operations to manage this information. The choice between counting elements, verifying data presence, or sampling

data subsets holds the key to significant performance and resource utilization impacts. This exploration is dedicated to unveiling the scenarios each Spark action best fits and guiding you through selecting the most efficient data handling techniques for your Apache Spark projects.

Use Cases

Delve into real-world applications for each Spark action:

- `count()` : Essential for when precision is non-negotiable, learn how this function serves as the backbone for reporting and decision-making based on dataset size.
- `isEmpty()` : Discover how this quick check can streamline your workflows, optimizing conditional logic and conserving valuable resources.
- `take(n)` : Uncover the advantages of inspecting subsets for efficient debugging, exploratory data analysis, or preliminary checks without the overhead of full dataset processing.

Functionality

`count():`

The `count()` method triggers a full scan of all data in the DataFrame, which means every row is evaluated. In distributed systems like Spark, this operation involves aggregating counts from various partitions and nodes, which can be time-consuming and resource-intensive. This method is deterministic, meaning it will return the same result each time for the same dataset.

Purpose: Returns the number of rows in a DataFrame.

When to use: Use `count()` when you need the exact number of records. It's particularly useful when you need to know the size of your dataset for reporting, logging, or conditional logic that depends on the size of the DataFrame.

Performance Considerations: Can be *expensive* in terms of computation, especially for large DataFrames, because Spark has to scan all the rows to calculate the count. If your DataFrame is distributed across many nodes, this operation requires shuffling and aggregation, which adds to the computational overhead.

Example: Picture yourself counting coins meticulously to ensure you have the exact amount you need. This time-consuming process is like using `count()` to go

through every row of your DataFrame, ensuring you know the total count but requiring time and attention to detail



isEmpty():

Unlike `count()`, `isEmpty()` is designed to be a quick check. Internally, Spark can achieve this by querying only a subset of the data or a single partition and immediately returning if any record is found, thus avoiding the full dataset scan that `count()` requires. This makes `isEmpty()` highly efficient, particularly for large, distributed datasets where a full scan would be costly.

Purpose: Checks whether a DataFrame is empty (i.e., has no rows).

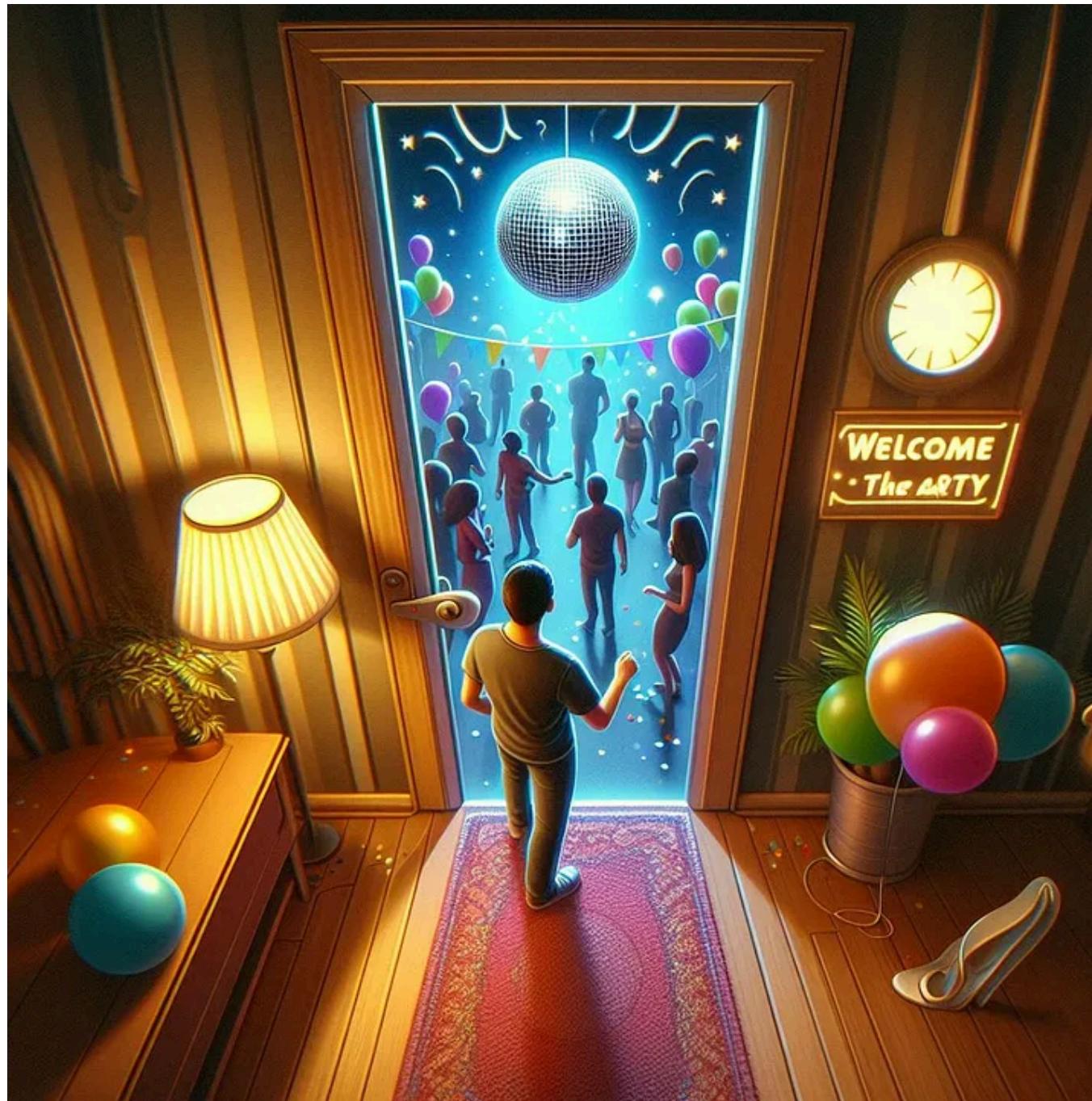
When to use: Use `isEmpty()` when you only need to know if the DataFrame contains any data and do not care about the exact count. It's useful for conditional logic where actions depend on the presence or absence of data.

Performance Considerations: Typically *faster than* `count()`, especially for large datasets, because it can stop as soon as it finds the first record. It's more efficient when you simply want to check the existence of any data rather than counting the total number of records.

Practical Applications: `isEmpty()` is particularly useful in control flow to prevent unnecessary operations on empty DataFrames, such as data

processing or writing to storage, which can save computational resources and time.

Example: Visualize walking up to a party and peeking inside the door to see if anyone's there. This quick glance, your "Immediate Entrance Check," parallels the `isEmpty()` function, offering a fast way to determine if there's any "life" or data in your DataFrame without fully entering or going through all the data.



take(n):

The `take(n)` method retrieves a specified number of records from the DataFrame. While similar to a SQL LIMIT clause, `take(n)` is executed in the context of the driver program, meaning the data is collected to the driver. This can be an expensive operation if n is large, as it involves network transfer of data from executors to the driver. However, for small n, it's a quick way to sample data.

Purpose: Returns the first n rows of the DataFrame.

When to use: Use `take()` when you need to retrieve a small number of rows for previewing, debugging, or testing. It allows you to quickly inspect a few records without scanning the entire DataFrame.

Performance Considerations: Usually *faster than* `count()` for large datasets because it only retrieves a specified number of rows. However, the performance depends on the value of n and the distribution of data across the cluster. For a small number of rows, it's efficient, but as n increases, the performance cost can also increase.

Example: Consider sampling chocolates from a box, picking just a few to taste.

This “Tasting Nibbles Neatly” action resembles the `take(n)` method, where you only look at a select number of rows to get a flavor of your data without needing to examine the entire “box” or DataFrame.



Notebook Execution Result

A delta table with 100 million records and executed the `count()`, `isEmpty()`, and `take(n)` actions on the dataframe.

```
from pyspark.sql import SparkSession

spark = SparkSession.builder \
    .appName("100_ml_table_read") \
    .getOrCreate()

delta_df = spark.sql("select * from table_action where city='City_9'")
```

```
delta_df: pyspark.sql.dataframe.DataFrame = [id: integer, name: string ... 3 more fields]
```

Retrieving dataframe where city='City_9'. Note- city is a partition column

```
1 print(delta_df.count())
```

▼ (3) Spark Jobs

- ▼ Job 0 [View](#) (Stages: 1/1)
 - Stage 0 32/32 succeeded [View](#)
- ▼ Job 1 [View](#) (Stages: 1/1, 1 skipped)
 - Stage 1 0/32 succeeded [View](#) skipped
 - Stage 2 1/1 succeeded [View](#)
- ▼ Job 2 [View](#) (Stages: 1/1, 2 skipped)
 - Stage 3 0/32 succeeded [View](#) skipped
 - Stage 4 0/1 succeeded [View](#) skipped
 - Stage 5 1/1 succeeded [View](#)

10000000

Command took 5.85 seconds

Usage of `count()` – Stages of Task

```
1 print(delta_df.isEmpty())
```

▼ (1) Spark Jobs

- ▼ Job 4 [View](#) (Stages: 1/1)
 - Stage 7 1/1 succeeded [View](#)

False

Command took 0.28 seconds

Usage of `isEmpty()`

```
1 print(delta_df.take(5))
```

► (1) Spark Jobs

- Job 6 [View](#) (Stages: 1/1)

```
[Row(id=1249279, name='Name_1249279', age=79, city='City_9', salary=1249278976.0), Row(id=1249269, name='Name_1249269', age=69, city='City_9', salary=1249268992.0), Row(id=1249259, name='Name_1249259', age=59, city='City_9', salary=1249259000.0), Row(id=1249249, name='Name_1249249', age=49, city='City_9', salary=1249249024.0), Row(id=1249239, name='Name_1249239', age=39, city='City_9', salary=1249239040.0)]
```

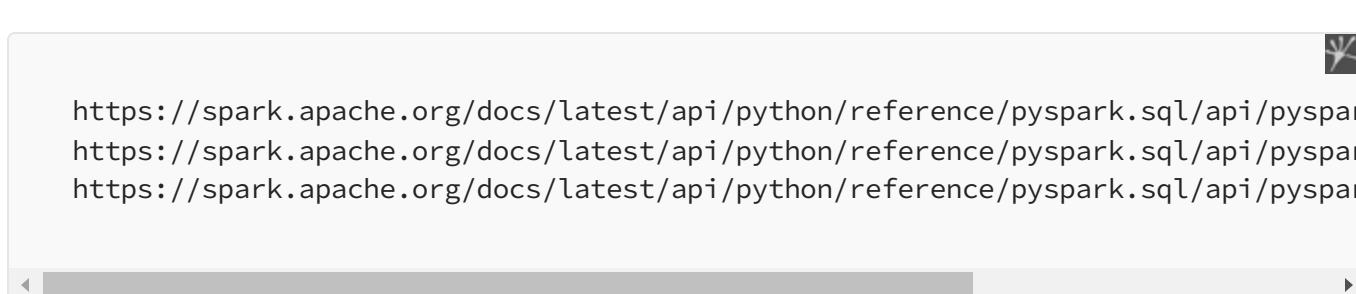
Command took 0.21 seconds

Usage of `take(n)`

Conclusion

In Spark, the choice between `count()`, `isEmpty()`, and `take(n)` should be made based on the specific requirements of the operation and the size of the dataset. While `count()` provides an exact record count, it does so at the cost of scanning the entire dataset. `isEmpty()` offers a performance-optimized alternative for simply checking the presence of data, making it suitable for conditional logic without the overhead of a full scan. `take(n)` is useful for quickly inspecting data but should be used judiciously to avoid excessive data transfer to the Spark driver. Understanding the internal workings and implications of these methods can help in writing more efficient Spark code and making better-informed decisions in data processing workflows.

Reference Links



The screenshot shows a browser window with three links listed vertically:

- <https://spark.apache.org/docs/latest/api/python/reference/pyspark.sql/api/pyspark.sql.Column.html#pyspark.sql.Column.count>
- <https://spark.apache.org/docs/latest/api/python/reference/pyspark.sql/api/pyspark.sql.Column.html#pyspark.sql.Column.isEmpty>
- <https://spark.apache.org/docs/latest/api/python/reference/pyspark.sql/api/pyspark.sql.Column.html#pyspark.sql.Column.take>

Thank you for taking the time to read this blog.

*If like this content FOLLOW & feel free to give this post a clap & comment
to help it reach wider audience.*

If you found it enjoyable or informative, please consider leaving a comment and sharing it with others.

Stay tuned for my upcoming blog posts !!!!!

Apache Spark

Spark

Data Engineering

Data Engineer

Pyspark



Written by Asish Panda

[Follow](#)

15 Followers

More from Asish Panda

 Asish Panda

Enhancing Lakehouse Table Design with Liquid Clustering

The flexible organizer for your data—adapt your layout on the fly without the need for a...

Aug 18  1

...

Apr 3  2

...

 Asish Panda

Small file problem in delta lake

Small Files, Big Impact: Enhancing Delta Lake Efficiency

Aug 4  22  1

 Asish Panda

Track and Audit Table Changes in Azure SQL Database: A...

In a vast digital library, AzureSQL tracked every book's movement, ensuring the...

Aug 19  

See all from Asish Panda

Recommended from Medium



 Abdur Rahman in Stackademic

Python is No More The King of Data Science

5 Reasons Why Python is Losing Its Crown

★ Oct 23 2.8K 19

 ...

 Hitesh Parab

Databricks Data Quality Framework using Great...

Author : Hitesh Parab & Yash Dholam

★ Sep 19 5

 ...

Lists



Natural Language Processing

1792 stories · 1405 saves



Staff Picks

756 stories · 1421 saves

Software Development Engineer Mar. 2020 – May 2021

- Developed Amazon checkout and payment services to handle traffic of 10 Million daily global transactions
- Integrated Iframes for credit cards and bank accounts to secure 80% of all consumer traffic and prevent CSRF, cross-site scripting, and cookie-jacking
- Led Your Transactions implementation for JavaScript front-end framework to showcase consumer transactions and reduce call center costs by \$25 Million
- Recovered Saudi Arabia checkout failure impacting 4000+ customers due to incorrect GET form redirection

Projects

NinjaPrep.io (React)

- Platform to offer coding problem practice with built in code editor and written + video solutions in React
- Utilized Nginx to reverse proxy IP address on Digital Ocean hosts
- Developed using Styled-Components for 95% CSS styling to ensure proper CSS scoping
- Implemented Docker with Seccomp to safely run user submitted code with < 2.2s runtime

HeatMap (JavaScript)

- Visualized Google Takeout location data of location history using Google Maps API and Google Maps heatmap code with React
- Included local file system storage to reliably handle 5mb of location history data
- Implemented Express to include routing between pages and jQuery to parse Google Map and implement heatmap overlay



 Alexander Nguyen in Level Up Coding

The resume that got a software engineer a \$300,000 job at Google.

1-page. Well-formatted.

★ Jun 1 25K 490



 F. Perry Wilson, MD MSCE 

How Old Is Your Body? Stand On One Leg and Find Out

 Mukovhe Mukwevho

Optimizing PySpark: Cutting Run-Times from 30 Minutes to Under ...

Recently, I encountered a challenge with a PySpark job that was taking over 30 minute...

★ Oct 7 219 3



 Anand Satheesh

Apache Spark Commonly seen errors in production and their...

According to new research, the time you can stand on one leg is the best marker of...

◆ Oct 23 9.5K 198

W+ ··· Jul 1 102

Apache Spark is a powerful tool for big data processing, it uses distributed data...

W+ ···

See more recommendations

Help Status About Careers Press Blog Privacy Terms Text to speech Teams