Spark by {Examples}  Spark (https://sparkbyexamples.com/)

## Spark Tutorial

Spark – Installation on Windows (https://sparkbyexamples.com/spark/apache-spark-installation-on-windows/)

Spark – Installation on Linux | Ubuntu (https://sparkbyexamples.com/spark/spark-installation-on-linux-ubuntu/)

Spark – Cluster Setup with Hadoop Yarn (https://sparkbyexamples.com/spark/spark-setup-on-hadoop-yarn/)

Spark – Web/Application UI (https://sparkbyexamples.com/spark/spark-web-ui-understanding/)

Spark – Setup with Scala and IntelliJ (https://sparkbyexamples.com/spark/spark-setup-run-with-scala-intellij/)

Spark – How to Run Examples From this Site on IntelliJ IDEA (https://sparkbyexamples.com/spark/how-to-run-spark-examples-from-intellij/)

Spark – SparkSession (https://sparkbyexamples.com/spark/sparksession-explained-with-examples/)

Spark – SparkContext (https://sparkbyexamples.com/spark/spark-sparkcontext/)

## Spark RDD Tutorial

Spark RDD – Parallelize (https://sparkbyexamples.com/apache-spark-rdd/how-to-create-an-rdd-using-parallelize/)

PySpark   (https://sparkbyexamples.com/pyspark-tutorial/)

Hive   (https://sparkbyexamples.com/apache-hive-tutorial/)

HBase   (https://sparkbyexamples.com/apache-hbase-tutorial/)

Kafka   (https://sparkbyexamples.com/apache-kafka-tutorials-with-examples/)

More   (https://sparkbyexamples.com/)

# Spark Performance Tuning & Best Practices

NNK (https://sparkbyexamples.com/author/admin/)  -  Apache Spark (https://sparkbyexamples.com/category/spark/)  /  PySpark (https://sparkbyexamples.com/category/pyspark/)

[ FAQ's ] (https://sparkbyexamples.com/spark-questions/)

Spark Performance tuning is a process to improve the performance of the Spark and PySpark applications by adjusting and optimizing system resources (CPU cores and memory), tuning some configurations, and following some framework guidelines and best practices.

Spark application performance can be improved in several ways. Below are the different articles I've written to cover these.

- Spark Guidelines and Best Practices (Covered in this article)
- Tuning System Resources (executors, CPU cores, memory) – In progress
- Tuning Spark Configurations (https://sparkbyexamples.com/spark/spark-sql-performance-tuning-configurations/) (AQE, Partitions e.t.c)

In this article, I have covered some of the framework guidelines and best practices to follow while developing Spark applications which ideally improves the performance of the application, most of these best practices would be the same for both Spark with Scala or PySpark (Python).

# Spark Performance Tuning – Best Guidelines & Practices

Spark performance tuning and optimization is a bigger topic which consists of several techniques, and configurations (resources memory & cores), here I've covered some of the best guidelines I've used to improve my workloads and I will keep updating this as I come acrossnew ways.

- Use DataFrame/Dataset over RDD
- Use coalesce() over repartition()
- Use mapPartitions() over map()
- Use Serialized data format's
- Avoid UDF's (User Defined Functions)
- Caching data in memory
  (https://sparkbyexamples.com/wp-
  admin/post.php?
  post=7347&action=edit#cache)
- Reduce expensive Shuffle operations
- Disable DEBUG & INFO Logging
  (https://sparkbyexamples.com/wp-
  admin/post.php?
  post=7347&action=edit#debug)

# 1. Use DataFrame/Dataset over RDD

For Spark jobs, prefer using Dataset/DataFrame over RDD as Dataset and DataFrame's includes several optimization modules to improve the performance of the Spark workloads. In PySpark use, DataFrame over RDD as Dataset's are not supported in PySpark applications.



Spark RDD is a building block of Spark programming, even when we use DataFrame/Dataset, Spark internally uses RDD to execute operations/queries but the efficient and optimized way by analyzing your query and creating the execution plan thanks to Project Tungsten (https://databricks.com/blog/2015/04/28/

project-tungsten-bringing-spark-closer-to-bare-metal.html) and Catalyst optimizer (https://databricks.com/glossary/catalyst-optimizer).

## Why RDD is slow?

Using RDD directly leads to performance issues as Spark doesn't know how to apply the optimization techniques and RDD serialize and de-serialize the data when it distributes across a cluster (repartition & shuffling).

Serialization and de-serialization are very expensive operations for Spark applications or any distributed systems, most of our time is spent only on serialization of data rather than executing the operations hence try to avoid using RDD.



Source: Databricks

## Is DataFrame Faster

Since Spark DataFrame maintains the structure of the data and column types (like an RDMS table) it can handle the data better by storing and managing more efficiently.

The DataFrame API does two things that help to do this (through the Tungsten project). First, using off-heap storage for data in binary format. Second, generating encoder code on the fly to work with this binary format for your specific objects.

Since Spark/PySpark DataFrame internally stores data in binary there is no need of Serialization and deserialization data when it distributes across a cluster hence you would see a performance improvement.

## Project Tungsten

Spark Dataset/DataFrame includes Project Tungsten (https://databricks.com/blog/2015/04/28/project-tungsten-bringing-spark-closer-to-bare-metal.html) which optimizes Spark jobs for **Memory and CPU efficiency**. Tungsten is a Spark SQL component that provides increased performance by rewriting Spark operations in bytecode, at runtime. Tungsten performance by focusing on jobs close to bare metal CPU and memory efficiency.

Since DataFrame is a column format that contains additional metadata, hence Spark can perform certain optimizations on a query. Before your query is run, a logical plan is created using Catalyst Optimizer (https://databricks.com/glossary/catalyst-optimizer) and then it's executed using the Tungsten execution engine.

## *What is Catalyst?*

Catalyst Optimizer
(https://databricks.com/glossary/catalyst-
optimizer) is an integrated query
optimizer and execution scheduler for
Spark Datasets/DataFrame. Catalyst
Optimizer is the place where Spark
tends to improve the speed of your code
execution by logically improving it.

Catalyst Optimizer can perform
refactoring complex queries and decides
the order of your query execution by
creating a rule-based and code-based
optimization.

Additionally, if you want type safety at
compile time prefer using Dataset. For
example, if you refer to a field that
doesn't exist in your code, Dataset
generates compile-time error whereas
DataFrame compiles fine but returns an
error during run-time.

## 2. Use coalesce() over repartition()

When you want to reduce the number of
partitions prefer using `coalesce()` as it
is an optimized or improved version
of `repartition()`
(https://sparkbyexamples.com/spark/spar
k-repartition-vs-coalesce/) where **the
movement of the data across the
partitions is lower using coalesce**
which ideally performs better when you
dealing with bigger datasets.

**Note:** Use repartition() when you wanted
to increase the number of partitions.

## Example repartition()

```
val rdd1 = spark.sparkContext.pa
  println("parallelize : "+rdd1.

val rdd2 = rdd1.repartition(4)
  println("Repartition size : "+
  rdd2.saveAsTextFile("/tmp/re-p
```

This yields output `Repartition size : 4` and the repartition re-distributes the data(as shown below) from all partitions which is full shuffle leading to very expensive operation when dealing with billions and trillions of data. By tuning the partition size to optimal, you can improve the performance of the Spark application

```
Partition 1 : 1 6 10 15 19
Partition 2 : 2 3 7 11 16
Partition 3 : 4 8 12 13 17
Partition 4 : 0 5 9 14 18
```

## Example coalesce()

```
val rdd3 = rdd1.coalesce(4)
  println("Repartition size : "+
  rdd3.saveAsTextFile("/tmp/coal
```

If you compared the below output with section 1, you will notice partition 3 has been moved to 2 and Partition 6 has moved to 5, resulting data movement from just 2 partitions.

```
Partition 1 : 0 1 2
Partition 2 : 3 4 5 6 7 8 9
Partition 4 : 10 11 12
Partition 5 : 13 14 15 16 17 18
```

# 3. Use mapPartitions() over map()

Spark `map()` and `mapPartitions()` transformation applies the function on each element/record/row of the DataFrame/Dataset and returns the new DataFrame/Dataset. mapPartitions() over map() prefovides performance improvement (https://sparkbyexamples.com/spark/spark-map-vs-mappartitions-transformation/) when you have havy initializations like initializing classes, database connections e.t.c

Spark `mapPartitions()` provides a facility to do heavy initializations (for example Database connection) once for each partition instead of doing it on every DataFrame row. This helps the performance of the Spark jobs when you dealing with heavy-weighted initialization on larger datasets.

## Example map()

```
import spark.implicits._
  val df3 = df2.map(row=>{
    val util = new Util() // Ini
    val fullName = util.combine(
    (fullName, row.getString(3),
  })
  val df3Map =  df3.toDF("fullNa
```

## Example mapPartitions()

**PySpark DataFrame**

```scala
val df4 = df2.mapPartitions(it
  val util = new Util()
  val res = iterator.map(row=>
    val fullName = util.combin
    (fullName, row.getString(3
  })
  res
})
val df4part = df4.toDF("fullNa
```

**Note:** One key point to remember is these both transformations returns the `Dataset[U]` but not the `DataFrame` (In Spark 2.0, `DataFrame = Dataset[Row]`).

## 4. Use Serialized data format's

Most of the Spark jobs run as a pipeline where one Spark job writes data into a File and another Spark jobs read the data, process it, and writes to another file for another Spark job to pick up. When you have such use case, prefer writing an intermediate file in Serialized and optimized formats like Avro, Kryo, Parquet e.t.c, any transformations on these formats performs better than text, CSV (https://sparkbyexamples.com/spark/spark-read-csv-file-into-dataframe/), and JSON (https://sparkbyexamples.com/spark/spark-read-and-write-json-file/).

## What is Parquet

Apache Parquet is a columnar file format that provides optimizations (https://sparkbyexamples.com/spark/spark-read-write-dataframe-parquet-example/) to speed up queries and is a

far more efficient file format than CSV or JSON, supported by many data processing systems.

It is compatible with most of the data processing frameworks in the Hadoop (https://en.wikipedia.org/wiki/Hadoop) echo systems. It provides efficient data compression and encoding schemes with enhanced performance to handle complex data in bulk.

val dF = spark.read.parquet("/tmp/output/people.parquet") //Read Parquet file df.write.parquet("/tmp/output/people-new.parquet")//Writing parquet file

## What is Avro

Apache Avro (https://sparkbyexamples.com/spark/read-write-avro-file-spark-dataframe/) is an open-source, row-based, data serialization and data exchange framework for Hadoop projects, originally developed by databricks as an open-source library that supports reading and writing data in Avro file format. it is mostly used in Apache Spark especially for Kafka-based data pipelines. When Avro data is stored in a file, its schema is stored with it, so that files may be processed later by any program.

It has build to serialize and exchange big data between different Hadoop based projects. It serializes data in a compact binary format and schema is in JSON format that defines the field names and data types.

```
val df = spark.read.format("avro
df.write.format("avro").save("pe
//Avro Spark SQL
spark.sqlContext.sql("CREATE TEM
    OPTIONS (path \"person.avro\
spark.sqlContext.sql("SELECT * F
```

## 5. Avoid UDF's (User Defined Functions)

Try to avoid Spark/PySpark UDF's (https://sparkbyexamples.com/spark/spark-sql-udf/) at any cost and use when existing Spark built-in functions are not available for use. UDF's are a black box to Spark hence it can't apply optimization and you will lose all the optimization Spark does on Dataframe/Dataset. When possible you should use Spark SQL built-in functions (https://sparkbyexamples.com/spark/spark-sql-functions-understanding/) as these functions provide optimization

Before you create any UDF, do your research to check if the similar function you wanted is already available in Spark SQL Functions (https://sparkbyexamples.com/spark/spark-sql-functions-understanding/). Spark SQL provides several predefined common functions and many more new functions are added with every release. hence, It is best to check before you reinventing the wheel.

## 6. Persisting & Caching data in memory

Spark persisting/caching is one of the best techniques to improve the performance of the Spark workloads. Spark *Cache* and P*ersist* are

optimization techniques in DataFrame / Dataset (https://sparkbyexamples.com/spark/spark-dataframe-cache-and-persist-explained/) for iterative and interactive Spark applications to improve the performance of Jobs.

Using `cache()` and `persist()` methods, Spark provides an optimization mechanism to store the intermediate computation of a Spark DataFrame so they can be reused in subsequent actions.

When you persist a dataset, each node stores it's partitioned data in memory and reuses them in other actions on that dataset. And Spark's persisted data on nodes are fault-tolerant meaning if any partition of a Dataset is lost, it will automatically be recomputed using the original transformations that created it.

```
df.where(col("State") === "PR").
```

When caching use in-memory columnar format, By tuning the batchSize property you can also improve Spark performance.

```
spark.conf.set("spark.sql.inMemo
spark.conf.set("spark.sql.inMemo
```

Spark provides several storage levels to store the cached data (https://sparkbyexamples.com/spark/spark-persistence-storage-levels/), use the once which suits your cluster.

# 7. Reduce expensive Shuffle operations

Shuffling is a mechanism Spark uses to redistribute the data (https://sparkbyexamples.com/spark/spark-repartition-vs-coalesce/) across different executors and even across machines. Spark shuffling triggers when we perform certain transformation operations like gropByKey(), reducebyKey(), join() on RDD and DataFrame.

Spark Shuffle is an expensive operation since it involves the following

- Disk I/O
- Involves data serialization and deserialization
- Network I/O

We cannot completely avoid shuffle operations in but when possible try to reduce the number of shuffle operations removed any unused operations.

Spark provides spark.sql.shuffle.partitions configurations to control the partitions of the shuffle, By tuning this property you can improve Spark performance.

```
park.conf.set("spark.sql.shuffl
qlContext.setConf("spark.sql.sh
```

## Disable DEBUG & INFO ogging

s is one of the simple ways to rove the performance of Spark Jobs and can be easily avoided by following good coding principles. During the development phase of Spark/PySpark application, we usually write debug/info

messages to console using `println()` and logging to a file using some logging framework (log4j);

These both methods results I/O operations hence cause performance issues when you run Spark jobs with greater workloads. Before promoting your jobs to production make sure you review your code and take care of the following.

**Remove or convert all println() statements to log4j info/debug.**

```
logger.debug("Debug logging
logger.info("Info logging me
```

**Disable DEBUG/INFO by enabling ERROR/WARN/FATAL logging**

If you are using log4j.properties use the following or use appropriate configuration based on your logging framework and configuration method (XML vs properties vs yaml)

```
log4j.rootLogger=warn, stdout
```

Personally I've seen this in my project where our team written 5 log statements in a map() (https://sparkbyexamples.com/spark/spark-map-vs-mappartitions-transformation/) transformation; When we are processing 2 million records which resulted 10 million I/O operations and caused my job running for hrs. After disabling DEBUG & INFO logging I've witnessed jobs running in few mins.

**Note:** Spark workloads are increasingly bottlenecked by CPU and memory use rather than I/O and network, but still

avoiding I/O operations are always a good practice.

# Spark Performance Tuning Conclusion

Spark with Scala or Python (pyspark) jobs run on huge dataset's, when not following good coding principles and optimization techniques you will pay the price with performance bottlenecks, by following the topics I've covered in this article you will achieve improvement programmatically however there are other ways to improve the performance and tuning Spark jobs (by config & increasing resources) which I will cover in my next article.

## References

- https://databricks.com/blog/2016/07/1
  4/a-tale-of-three-apache-spark-apis-
  rdds-dataframes-and-datasets.html
  (https://databricks.com/blog/2016/07/1
  4/a-tale-of-three-apache-spark-apis-
  rdds-dataframes-and-datasets.html)
- https://databricks.com/blog/2015/04/2
  8/project-tungsten-bringing-spark-
  closer-to-bare-metal.html
  (https://databricks.com/blog/2015/04/2
  8/project-tungsten-bringing-spark-
  closer-to-bare-metal.html)

Hope you like this article, leave me a comment if you like it or have any questions.

Happy Learning !!

**Share this:**

 (https://sparkbyexamples.com/spark/spark-
performance-tuning/?share=facebook&nb=1)

 (https://sparkbyexamples.com/spark/spark-
performance-tuning/?share=reddit&nb=1)

PySpark – countDistinct()
(https://sparkbyexamples.com/py
spark/pyspark-count-distinct-
from-dataframe/)

PySpark – sum(), avg()
(https://sparkbyexamples.com/py
spark/pyspark-dataframe-
groupby-and-sort-by-descending-
order/)

PySpark – row_number()
(https://sparkbyexamples.com/py
spark/pyspark-window-
functions/#row_number)

PySpark – rank()
(https://sparkbyexamples.com/py
spark/pyspark-window-
functions/#rank)

PySpark – dense_rank()
(https://sparkbyexamples.com/py
spark/pyspark-window-
functions/#dense_rank)

PySpark – percent_rank()
(https://sparkbyexamples.com/py
spark/pyspark-window-
functions/#percent_rank)

PySpark – typedLit()
(https://sparkbyexamples.com/py
spark/pyspark-lit-add-literal-
constant/#typedlit)

PySpark – from_json()
(https://sparkbyexamples.com/py
spark/pyspark-json-functions-
with-examples/#from_json)

PySpark – to_json()
(https://sparkbyexamples.com/py
spark/pyspark-json-functions-
with-examples/#to_json)

PySpark – json_tuple()
(https://sparkbyexamples.com/py
spark/pyspark-json-functions-
with-examples/#json_tuple)

PySpark – get_json_object()
(https://sparkbyexamples.com/py
spark/pyspark-json-functions-
with-examples/#get_json_object)

PySpark – schema_of_json()
(https://sparkbyexamples.com/py

TAGS: **PYSPARK OPTIMIZATION
(HTTPS://SPARKBYEXAMPLES.COM/TAG/PYSPARK-
OPTIMIZATION/)**, **PYSPARK PERFORMANCE
(HTTPS://SPARKBYEXAMPLES.COM/TAG/PYSPARK-
PERFORMANCE/)**, **PYSPARK TUNING
(HTTPS://SPARKBYEXAMPLES.COM/TAG/PYSPARK-
TUNING/)**, **SPARK OPTIMIZATION
(HTTPS://SPARKBYEXAMPLES.COM/TAG/SPARK-
OPTIMIZATION/)**, **SPARK PERFORMANCE
(HTTPS://SPARKBYEXAMPLES.COM/TAG/SPARK-
PERFORMANCE/)**

## NNK
## (Https://Sparkbyexamples.Com/Author/Admin/)

(https://sp
arkbyexa
mples.co
m/author/
admin/)

SparkByExamples.com is a Big Data and Spark
examples community page, all examples are simple and
easy to understand and well tested in our development
environment Read more ..
(https://sparkbyexamples.com/about-sparkbyexamples/)

**❯ THIS POST HAS 4 COMMENTS**

### Sudhindra Kulkarni

**23 JUN 2021**          **REPLY**

Very nice explanation with
good examples. Please
keep the articles moving.

**NNK**    **6 JUL 2021**    **REPLY**

Thanks Sudhindra.

**Naveen**  **26 MAR 2021**  **REPLY**

Hi..
Please Post the Performance tuning the spark code to load oracle table..

**Anonymous**

**13 JAN 2021**  **REPLY**

Very Well Explained in Simple way

## Leave a Reply

Enter your comment here...

---

**Categories**

Apache Hadoop (https://sparkbyexamples.com/category/hadoop/)

Apache Spark (https://sparkbyexamples.com/category/spark/)

Apache Spark Streaming (https://sparkbyexamples.com/category/spark/apache-spark-streaming/)

**Recent Posts**

Spark regexp_replace() – Replace String Value (https://sparkbyexamples.com/spark/spark-regexp_replace-replace-string-value/)

How to Run a PySpark Script from Python? (https://sparkbyexamples.com/pyspark/run-pyspark-script-from-python-subprocess/)

**About SparkByExamples.Com**

SparkByExamples.com is a Big Data and Spark examples community page, all examples are simple and easy to understand, and well tested in our development environment Read more .. (https://sparkbyexamples.com/about-sparkbyexamples/)

**Follow Us**

(https: (https:

//www. //www.

(https: facebo linkedi (https:

//twitte ok.co n.com/ //githu

r.com/ m/spar in/n- b.com/

sparkb kbyex nk- spark-

yexam ample b860a examp

ples) s/) 8193/) les/)