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

Spark + (https://sparkbyexamples.com/)

## Spark Tutorial

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

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

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

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

### Top B.Tech Program in B'lore

B.Tech @ PES University, B'lore - Top Infrastructure | Quality Faculty | Best Placements

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

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

# Spark DataFrame Cache and Persist Explained

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

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

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

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

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

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

Spark *Cache* and *Persist* are optimization techniques in DataFrame / Dataset for iterative and interactive Spark applications to improve the performance of Jobs. In this article, you will learn What is Spark `cache()` and `persist()`, how to use it in DataFrame, understanding the difference between `Caching` and `Persistance` (https://sparkbyexamples.com/spark/spark-difference-between-cache-and-persist/) and how to use these two with DataFrame, and Dataset using Scala examples.

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

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

More ⌄ (https://sparkbyexamples.com/) 🔍

## Spark RDD Tutorial

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

Spark RDD – Read text file (https://sparkbyexamples.com/apache-spark-rdd/spark-read-multiple-text-files-into-a-single-rdd/)

Though Spark provides computation 100 x times faster than traditional Map Reduce jobs, If you have not designed the jobs to reuse the repeating computations you will see degrade in performance when you are dealing with billions or trillions of data. Hence, we may need to look at the stages and use optimization techniques as one of the ways to improve performance.

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.

# Advantages for Caching and Persistence of DataFrame

Below are the advantages of using Spark Cache and Persist methods.

**Cost efficient** – Spark computations are very expensive hence reusing the computations are used to save cost.

**Time efficient** – Reusing the repeated computations saves lots of time.

**Execution time** – Saves execution time of the job and we can perform more jobs on the same cluster.

# Spark Cache Syntax and Example

Spark DataFrame or Dataset `cache()` method by default saves it to storage level `MEMORY_AND_DISK` because recomputing the in-memory columnar representation of the underlying table is expensive. Note that this is different from the default cache level of `RDD.cache()` which is '`MEMORY_ONLY`'.

**Syntax**

```
cache() : Dataset.this.type
```

Spark `cache()` method in Dataset class internally calls `persist()` method which in turn uses `sparkSession.sharedState.cacheManager.cacheQuery` to cache the result set of DataFrame or Dataset. Let's look at an example.

**Example**

```
val spark:SparkSession = Spar
  .master("local[1]")
  .appName("SparkByExamples.c
  .getOrCreate()

//read csv with options
val df = spark.read.options(M
  .csv("src/main/resources/zi

val df2 = df.where(col("State
df2.show(false)

println(df2.count())

val df3 = df2.where(col("Zipc

println(df2.count())
```

## DataFrame Persist Syntax and Example

Spark persist() method is used to store the DataFrame or Dataset to one of the storage levels `MEMORY_ONLY`,`MEMORY_AND_DISK`, `MEMORY_ONLY_SER`, `MEMORY_AND_DISK_SER`, `DISK_ONLY`, `MEMORY_ONLY_2`,`MEMORY_AND_DISK_2` and more.

Caching or persisting of Spark DataFrame or Dataset is a lazy operation, meaning a DataFrame will not be cached until you trigger an action.

**Syntax**

```
1) persist() : Dataset.this.typ
2) persist(newLevel : org.apach
```

Spark persist has two signature first signature doesn't take any argument which by default saves it to `MEMORY_AND_DISK` storage level and the second signature which takes `StorageLevel` as an argument to store it to different storage levels.

**Example**

```
val dfPersist = df.persist()
dfPersist.show(false)
```

Using the second signature you can save DataFrame/Dataset to any storage levels.

```
val dfPersist = df.persist(St
dfPersist.show(false)
```

This stores DataFrame/Dataset into Memory.

Note that Dataset `cache()` is an alias for `persist(StorageLevel.MEMORY_AND_DISK)`

# Unpersist syntax and Example

Spark automatically monitors every persist() and cache() calls you make and it checks usage on each node and drops persisted data if not used or by using least-recently-used (LRU) algorithm. You can also manually remove using `unpersist()` method. unpersist() marks the Dataset as non-persistent, and remove all blocks for it from memory and disk.

**Syntax**

```
unpersist() : Dataset.this.type
unpersist(blocking : scala.Bool
```

E**xample**

```
val dfPersist = dfPersist.unp
```

unpersist(Boolean) with boolean as
argument blocks until all blocks are
deleted.

## Spark Persist storage levels

All different storage level Spark
supports are available at
`org.apache.spark.storage.Storag
eLevel` class. The storage level
specifies how and where to persist or
cache a Spark DataFrame and Dataset.

`MEMORY_ONLY` – This is the default
behavior of the RDD `cache()` method
and stores the RDD or DataFrame as
deserialized objects to JVM memory.
When there is no enough memory
available it will not save DataFrame of
some partitions and these will be re-
computed as and when required. This
takes more memory. but unlike RDD,
this would be slower than
`MEMORY_AND_DISK` level as it
recomputes the unsaved partitions and
recomputing the in-memory columnar
representation of the underlying table is
expensive

`MEMORY_ONLY_SER` – This is the same
as `MEMORY_ONLY` but the difference
being it stores RDD as serialized
objects to JVM memory. It takes lesser
memory (space-efficient) then
MEMORY_ONLY as it saves objects as
serialized and takes an additional few
more CPU cycles in order to
deserialize.

`MEMORY_ONLY_2` – Same as `MEMORY_ONLY` storage level but replicate each partition to two cluster nodes.

`MEMORY_ONLY_SER_2` – Same as `MEMORY_ONLY_SER` storage level but replicate each partition to two cluster nodes.

`MEMORY_AND_DISK` – This is the default behavior of the DataFrame or Dataset. In this Storage Level, The DataFrame will be stored in JVM memory as a deserialized objects. When required storage is greater than available memory, it stores some of the excess partitions into disk and reads the data from disk when it required. It is slower as there is I/O involved.

`MEMORY_AND_DISK_SER` – This is same as `MEMORY_AND_DISK` storage level difference being it serializes the DataFrame objects in memory and on disk when space not available.

`MEMORY_AND_DISK_2` – Same as `MEMORY_AND_DISK` storage level but replicate each partition to two cluster nodes.

`MEMORY_AND_DISK_SER_2` – Same as `MEMORY_AND_DISK_SER` storage level but replicate each partition to two cluster nodes.

`DISK_ONLY` – In this storage level, DataFrame is stored only on disk and the CPU computation time is high as I/O involved.

`K_ONLY_2` – Same `DISK_ONLY` storage level but icate each partition to two cluster es.

## onclusion

his article, you have learned Spark he() and `persist()` methods are used as optimization techniques to save interim computation results of DataFrame or Dataset and reuse them

sequently and learned what is the
erence between Spark Cache and
sist and finally saw their syntaxes
usages with Scala examples.

py Learning !!

## ference

ttps://spark.apache.org/docs/latest/r
d-programming-guide.html#rdd-
ersistence
https://spark.apache.org/docs/latest/
dd-programming-guide.html#rdd-
ersistence)

---

e this:

(https://sparkbyexamples.com/spark/spark-
aframe-cache-and-persist-explained/?
re=facebook&nb=1)

(https://sparkbyexamples.com/spark/spark-
aframe-cache-and-persist-explained/?
re=reddit&nb=1)

(https://sparkbyexamples.com/spark/spark-
aframe-cache-and-persist-explained/?
re=pinterest&nb=1)

(https://sparkbyexamples.com/spark/spark-
dataframe-cache-and-persist-explained/?
share=tumblr&nb=1)

(https://sparkbyexamples.com/spark/spark-
dataframe-cache-and-persist-explained/?
share=pocket&nb=1)

(https://sparkbyexamples.com/spark/spark-
dataframe-cache-and-persist-explained/?
share=linkedin&nb=1)

(https://sparkbyexamples.com/spark/spark-
dataframe-cache-and-persist-explained/?
share=twitter&nb=1)

**TAGS:** **CACHE (HTTPS://SPARKBYEXAMPLES.COM/TAG/CACHE/)**, **OPTIMIZATION (HTTPS://SPARKBYEXAMPLES.COM/TAG/OPTIMIZATION/)**, **PERSIST (HTTPS://SPARKBYEXAMPLES.COM/TAG/PERSIST/)**.

## NNK (Https://Sparkbyexamples.Com/Author/Admin/)
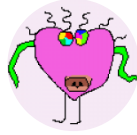
(https://sparkbyexamples.com/author/admin/)

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/)

**Arun**  24 MAR 2021  REPLY

Please update color of code snippets to dark shades current ones are not clearly visible.

**NNK** 24 MAR 2021  REPLY

Hi Arun, I don't understand, right now I have a black background and code in white/gray. Could you please let me know what is not clearly visible?

**Anonymous**

15 JUN 2020  REPLY

please remove the autoscroll. Its not helping.

**NNK** 15 JUN 2020  REPLY

Could you please let me know what browser you are using? I am not seeing auto scroll on Chrome?

## Leave a Reply

ark-check-string-column-has-numeric-
values/)

Spark Check Column Data Type is
Integer or String
(https://sparkbyexamples.com/spark/sp
ark-check-column-data-type-is-integer-
or-string/)

ark-check-string-column-has-numeric-
values/)

Spark Check Column Data Type is
Integer or String
(https://sparkbyexamples.com/spark/sp
ark-check-column-data-type-is-integer-
or-string/)