# 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/)

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/)

### Advance Certification Program

Learn from IIT Madras Professors & Industry Practitioners. Advance Your Career Now!

# Spark Read and Write Apache Parquet

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


(https://sparkbyexamples.com/spark-and-parquet-example/spark-parquet-example/)

Example of Spark read & write parquet file

In this tutorial, we will learn what is Apache Parquet?, It's advantages and how to read from and write Spark DataFrame to Parquet file format using Scala example. The example provided here is also available at Github

(https://github.com/sparkbyexamples/spark-examples/blob/master/spark-sql-examples/src/main/scala/com/sparkbyexamples/spark/dataframe/ParquetExample.scala) repository for reference.



- Apache Parquet Introduction
- Apache Parquet Advantages
- Spark Write DataFrame to Parquet file format
- Spark Read Parquet file into DataFrame
- Appending to existing Parquet file
- Running SQL queries
- Partitioning and Performance Improvement
- Reading a specific Parquet Partition
- Spark parquet schema

# Apache Parquet

# Introduction

Apache Parquet (https://parquet.apache.org/) is a columnar file format that provides optimizations 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.

Spark SQL provides support for both reading and writing Parquet files that automatically capture the schema of the original data, It also reduces data storage by 75% on average. Below are some advantages of storing data in a parquet format. Spark by default supports Parquet in its library hence we don't need to add any dependency libraries.

# Apache Parquet Advantages:

Below are some of the advantages of using Apache Parquet. combining these benefits with Spark improves performance and gives the ability to work with structure files.



- **Reduces IO operations.**
- **Fetches specific columns that you need to access.**
- **It consumes less space.**
- **Support type-specific encoding.**

# Apache Parquet Spark Example

Before we go over the Apache parquet with the Spark example, first, let's Create a Spark DataFrame (https://sparkbyexamples.com/spark/different-ways-to-create-a-spark-dataframe/) from `Seq` object. Note that `toDF()` (https://sparkbyexamples.com/spark/different-ways-to-create-a-spark-dataframe/#from-collection) function on sequence object is available only when you import implicits using `spark.sqlContext.implicits._`. This complete spark parquet example is available at Github (https://github.com/sparkbyexamples/spark-examples/blob/master/spark-sql-examples/src/main/scala/com/sparkbyexamples/spark/dataframe/ParquetExample.scala) repository for reference.

```
val data = Seq(("James ","","Smi
                ("Michael ","Rose"
                ("Robert ","","Wil
                ("Maria ","Anne","
                ("Jen","Mary","Bro

val columns = Seq("firstname","m

import spark.sqlContext.implicit
val df = data.toDF(columns:_*)
```

The above example creates a data frame with columns "firstname", "middlename", "lastname", "dob", "gender", "salary"

## Spark Write DataFrame to Parquet file format

Using `parquet()` function of `DataFrameWriter` class, we can write Spark DataFrame to the Parquet file. As mentioned earlier Spark doesn't need any additional packages or libraries to use Parquet as it by default provides with Spark. easy isn't it? so we don't

have to worry about version and compatibility issues. In this example, we are writing DataFrame to "people.parquet" file.

```
df.write.parquet("/tmp/output/pe
```

Writing Spark DataFrame to Parquet format preserves the column names and data types, and all columns are automatically converted to be nullable for compatibility reasons. Notice that all part files Spark creates has parquet extension.



## Spark Read Parquet file into DataFrame

Similar to write, DataFrameReader provides parquet() function (spark.read.parquet) to read the parquet files and creates a Spark DataFrame. In this example snippet, we are reading data from an apache parquet file we have written before.

```
val parqDF = spark.read.parquet(
```

printing schema of DataFrame returns columns with the same names and data types.

# Append to existing Parquet file

Spark provides the capability to append DataFrame to existing parquet files using "append" save mode. In case, if you want to overwrite use "overwrite" save mode.

```
df.write.mode('append').parquet(
```

# Using SQL queries on Parquet

We can also create a temporary view on Parquet files and then use it in Spark SQL statements. This temporary table would be available until the SparkContext present.

```
parqDF.createOrReplaceTempView("
val parkSQL = spark.sql("select
```

Above predicate on spark parquet file does the file scan which is performance bottleneck like table scan on a traditional database. We should use partitioning in order to improve performance.
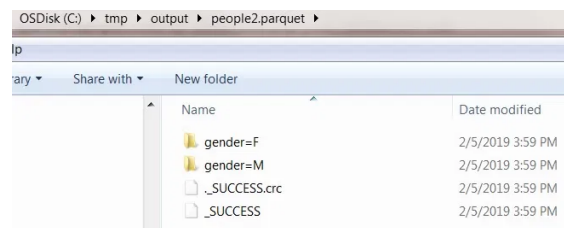
# Spark parquet partition – Improving performance

Partitioning is a feature of many databases and data processing frameworks and it is key to make jobs work at scale. We can do a parquet file partition using spark `partitionBy()` function.

```
df.write.partitionBy("gender","s
        .parquet("/tmp/output/pe
```

Parquet Partition creates a folder hierarchy for each spark partition; we have mentioned the first partition as gender followed by salary hence, it creates a salary folder inside the gender folder.



This is an example of how to write a Spark DataFrame by preserving the partitioning on gender and salary columns.

```
val parqDF = spark.read.parquet(
parqDF.createOrReplaceTempView("
val df = spark.sql("select * fro
```

The execution of this query is significantly faster than the query without partition (https://sparkbyexamples.com/spark/spark-performance-tuning/). It filters the data first on gender and then applies filters on salary.

# Spark Read a specific Parquet partition

```
val parqDF = spark.read.parquet(
```
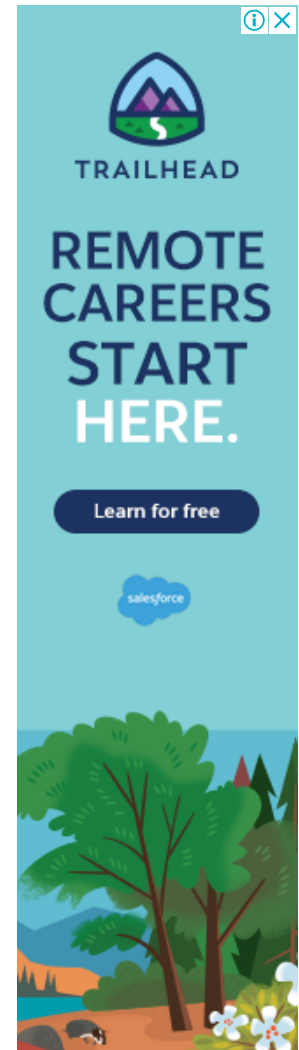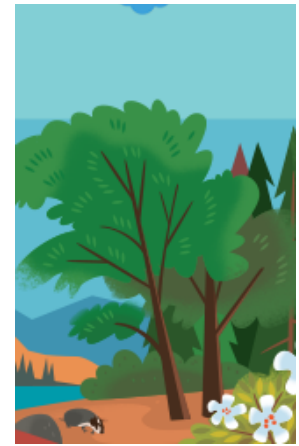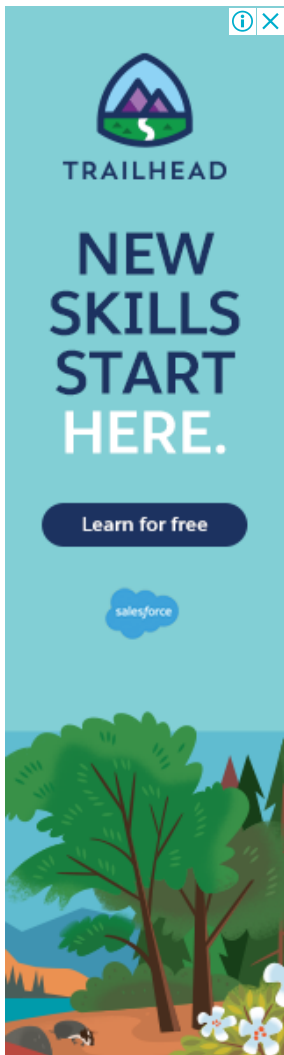
This code snippet retrieves the data from the gender partition value "M".

The complete code can be downloaded from GitHub (https://github.com/sparkbyexamples/spark-examples/blob/master/spark-sql-examples/src/main/scala/com/sparkbyexamples/spark/dataframe/ParquetExample.scala)

# Complete Spark Parquet Example

```scala
package com.sparkbyexamples.spar

import org.apache.spark.sql.Spar

object ParquetExample {

  def main(args:Array[String]):U

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

    val data = Seq(("James ","",
                   ("Michael ","Ro
                   ("Robert ","","
                   ("Maria ","Anne
                   ("Jen","Mary","

    val columns = Seq("firstname
    import spark.sqlContext.impl
    val df = data.toDF(columns:_
    df.show()
    df.printSchema()
    df.write
      .parquet("/tmp/output/peop
    val parqDF = spark.read.parq
    parqDF.createOrReplaceTempVi
    spark.sql("select * from Par
    val parkSQL = spark.sql("sel
    parkSQL.show()
    parkSQL.printSchema()
    df.write
      .partitionBy("gender","sal
      .parquet("/tmp/output/peop
    val parqDF2 = spark.read.par
    parqDF2.createOrReplaceTempV
    val df3 = spark.sql("select
    df3.explain()
    df3.printSchema()
    df3.show()
    val parqDF3 = spark.read
      .parquet("/tmp/output/peop
    parqDF3.show()
  }
}
```
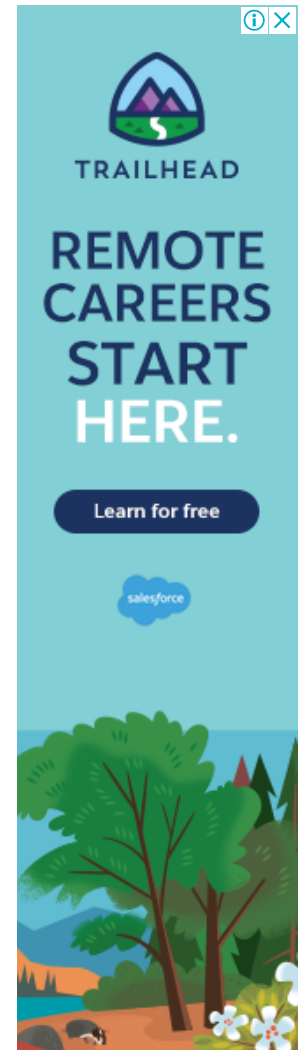
# Conclusion:

You have learned how to read a write an apache parquet data files in Spark and also learned how to improve the performance (https://sparkbyexamples.com/spark/spark-performance-tuning/) by using partition and filtering data with a partition key and finally appending to and overwriting existing parquet files.

Happy Learning !! 🙂

---

TAGS: **APACHE PARQUET (HTTPS://SPARKBYEXAMPLES.COM/TAG/APACHE-PARQUET/)**, **APACHE PARQUET SPARK (HTTPS://SPARKBYEXAMPLES.COM/TAG/APACHE-PARQUET-SPARK/)**, **SPARK READ PARQUET (HTTPS://SPARKBYEXAMPLES.COM/TAG/SPARK-READ-PARQUET/)**, **SPARK WRITE PARQUET (HTTPS://SPARKBYEXAMPLES.COM/TAG/SPARK-WRITE-PARQUET/)**

---

# Conclusion:

## 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/)

---

> **THIS POST HAS 5 COMMENTS**

### Anonymous

**3 MAR 2021**   **REPLY**

Great Work.

---

### NNK   **3 MAR 2021**   **REPLY**

Glad you like it.

---

### Kass   **2 JAN 2021**   **REPLY**

Nice work, thanks! Perhaps good to mention that partitioning is supported in various data formats (csv, json etc) and not just parquet.

---

### NNK   **2 JAN 2021**   **REPLY**

Hi Kass, Thanks for your feedback. Certainly, I will update the article with your suggestion.

---

### Pickard   **14 NOV 2019**   **REPLY**

Great tutorial, thank you!

## Leave a Reply

# Data Science Online Course

Live Instructor-Led Classes & Gain Hands-on Exposure to Data Science, R, SAS, Python & AI.

**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**

Apache Kafka
(https://sparkbyexamples.com/categor
y/kafka/)

Apache HBase
(https://sparkbyexamples.com/categor
y/hbase/)

Apache Cassandra
(https://sparkbyexamples.com/categor
y/cassandra/)

Snowflake Database
(https://sparkbyexamples.com/categor
y/snowflake/)

H2O Sparkling Water
(https://sparkbyexamples.com/categor
y/h2o-sparkling-water/)

PySpark
(https://sparkbyexamples.com/categor
y/pyspark/)

Spark SQL like() Using Wildcard
Example
(https://sparkbyexamples.com/spar
k-sql-like-using-wildcard-example/)

Spark isin() & IS NOT IN Operator
Example
(https://sparkbyexamples.com/spar
k-isin-is-not-in-operator-example/)

Spark – Get Size/Length of Array & Map
Column
(https://sparkbyexamples.com/spar
k-get-size-length-of-array-map-column/)

Spark Using Length/Size Of a
DataFrame Column
(https://sparkbyexamples.com/spar
k-using-length-size-of-a-dataframe-
column/)

Spark rlike() Working with Regex
Matching Examples
(https://sparkbyexamples.com/spar
k-rlike-regex-matching-examples/)

Spark Check String Column Has
Numeric Values
(https://sparkbyexamples.com/spar
k-check-string-column-has-numeric-
values/)

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

(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/)