Spark by {Examples} (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/)

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

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/)

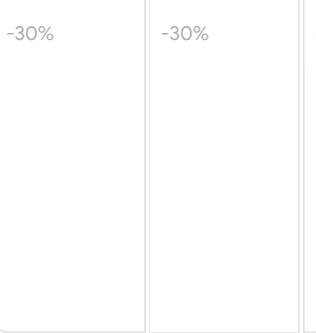FAQ's (https://sparkbyexamples.com/spark-questions/)

More ⌄ (https://sparkbyexamples.com/)

# Spark Merge Two DataFrames with Different Columns or Schema

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

In Spark or PySpark let's see how to merge/unione two DataFrames with a different number of columns (different schema). In Spark 3.1, you can easily achieve this using unionByName() transformation by passing allowMissingColumns with the value true. In order version, this property is not available

```
//Scala
merged_df = df1.unionByName(df2

#PySpark
merged_df = df1.unionByName(df2
```

The difference between
`unionByName()` function and `union()`
is that this function
resolves columns by name (not by
position). In other words,
unionByName() is used to merge two
DataFrame's by column names instead
of by position.

In case if you are using older than
Spark 3.1 version, use below approach
to merge DataFrame's with different
column names.

- Spark Merge DataFrames with
  Different Columns (Scala Example)
- PySpark Merge DataFrames with
  Different Columns (Python Example)

# Spark Merge Two DataFrames with Different Columns

In this section I will cover Spark with
Scala example of how to merge two
different DataFrames, first let's create
DataFrames with different number of
columns. DataFrame `df1` missing
column `state` and `salary` and `df2`
missing column `age`.

```
//Create DataFrame df1 with col
val data = Seq(("James","Sales"
                ("Robert","Sales
import spark.implicits._
val df1 = data.toDF("name","dep
df1.printSchema()

//root
// |-- name: string (nullable =
// |-- dept: string (nullable =
// |-- age: long (nullable = tr
```

Second DataFrame



```
//Create DataFrame df1 with col
val data2=Seq(("James","Sales",
                ("Jen","Finance",
val df2 = data2.toDF("name","de
df2.printSchema()

//root
// |-- name: string (nullable =
// |-- dept: string (nullable =
// |-- state: string (nullable
// |-- salary: long (nullable =
```

Now create a new DataFrames from existing after adding missing columns. newly added columns contains null values and we add constant column using lit() function (https://sparkbyexamples.com/spark/using-lit-and-typedlit-to-add-a-literal-or-constant-to-spark-dataframe/).

```
val merged_cols = df1.columns.t
import org.apache.spark.sql.fun
def getNewColumns(column: Set[S
    merged_cols.toList.map(x =>
      case x if column.contains
      case _ => lit(null).as(x)
    })
}
val new_df1=df1.select(getNewCo
val new_df2=df2.select(getNewCo
```

Finally merge two DataFrame's by using
column names

```
//Finally join two dataframe's
val merged_df=new_df1.unionByNa
merged_df.show()

//+-------+---------+----+-----
//|   name|     dept| age|state
//+-------+---------+----+-----
//|  James|    Sales|  34| null
//|Michael|    Sales|  56| null
//| Robert|    Sales|  30| null
//|  Maria|  Finance|  24| null
//|  James|    Sales|null|   NY
//|  Maria|  Finance|null|   CA
//|    Jen|  Finance|null|   NY
//|   Jeff|Marketing|null|   CA
//+-------+---------+----+-----
```

## PySpark Merge Two
## DataFrames with
## Different Columns

In PySpark to merge two DataFrames
with different columns, will use the
similar approach explain above and
uses `unionByName()` transformation.
First let's create DataFrame's with
different number of columns.

## Spark SQL Functions

```python
from pyspark.sql import SparkSe
spark = SparkSession.builder.ap

#Create DataFrame df1 with colu
data = [("James","Sales",34), (
    ("Robert","Sales",30), ("Ma
columns= ["name","dept","age"]
df1 = spark.createDataFrame(dat
df1.printSchema()

#Create DataFrame df1 with colu
data2=[("James","Sales","NY",90
    ("Jen","Finance","NY",7900)
columns2= ["name","dept","state
df2 = spark.createDataFrame(dat
df2.printSchema()
```

Now add missing columns `state` and `salary` to df1 and `age` to df2 with null values.

```python
#Add missing columns 'state' &
from pyspark.sql.functions impo
for column in [column for colum
    df1 = df1.withColumn(column

#Add missing column 'age' to df
for column in [column for colum
    df2 = df2.withColumn(column
```

Now merge/union the DataFrames using `unionByName()`. The difference between `unionByName()` function and `union()` is that this function resolves columns by name (not by position). In other words, unionByName() is used to merge two DataFrame's by column names instead of by position.

```
#Finally join two dataframe's d
merged_df=df1.unionByName(df2)
merged_df.show()
```

# Conclusion

In this article, you have learned with spark & PySpark examples of how to merge two DataFrames with different columns can be done by adding missing columns to the DataFrame's and finally union them using unionByName().

Happy Learning !!

**Share this:**

(https://sparkbyexamples.com/spark/spark-merge-two-dataframes-with-different-columns/?share=facebook&nb=1)

(https://sparkbyexamples.com/spark/spark-merge-two-dataframes-with-different-columns/?share=reddit&nb=1)

(https://sparkbyexamples.com/spark/spark-merge-two-dataframes-with-different-columns/?share=pinterest&nb=1)

(https://sparkbyexamples.com/spark/spark-merge-two-dataframes-with-different-columns/?share=tumblr&nb=1)

(https://sparkbyexamples.com/spark/spark-merge-two-dataframes-with-different-columns/?share=pocket&nb=1)

(https://sparkbyexamples.com/spark/spark-merge-two-dataframes-with-different-columns/?share=linkedin&nb=1)

(https://sparkbyexamples.com/spark/spark-merge-two-dataframes-with-different-columns/?share=twitter&nb=1)

## 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/)

# Leave a Reply

## PySpark DataFrame

## PySpark SQL Functions

## PySpark Datasources

## PySpark Built-In Functions

PySpark – when()
(https://sparkbyexamples.com/py
spark/pyspark-when-otherwise/)

PySpark – expr()
(https://sparkbyexamples.com/py
spark/pyspark-sql-expr-
expression-function/)

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

PySpark – split()
(https://sparkbyexamples.com/py
spark/pyspark-convert-string-to-
array-column/)

PySpark – concat_ws()
(https://sparkbyexamples.com/py
spark/pyspark-convert-array-
column-to-string-column/)

Pyspark – substring()
(https://sparkbyexamples.com/py
spark/pyspark-substring-from-a-
column/)

PySpark – translate()
(https://sparkbyexamples.com/py
spark/pyspark-replace-column-

## Categories

Apache Hadoop
(https://sparkbyexamples.com/catego
ry/hadoop/)

Apache Spark
(https://sparkbyexamples.com/catego
ry/spark/)

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

Apache Kafka
(https://sparkbyexamples.com/catego
ry/kafka/)

Apache HBase
(https://sparkbyexamples.com/catego
ry/hbase/)

Apache Cassandra
(https://sparkbyexamples.com/catego
ry/cassandra/)

Snowflake Database
(https://sparkbyexamples.com/catego
ry/snowflake/)

H2O Sparkling Water
(https://sparkbyexamples.com/catego
ry/h2o-sparkling-water/)

PySpark
(https://sparkbyexamples.com/catego
ry/pyspark/)

## Recent Posts

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

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

Spark SQL like() Using Wildcard
Example
(https://sparkbyexamples.com/spark/sp
ark-sql-like-using-wildcard-example/)

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

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

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

Spark rlike() Working with Regex
Matching Examples
(https://sparkbyexamples.com/spark/sp
ark-rlike-regex-matching-examples/)

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

Spark Check Column Data Type is
Integer or String
(https://sparkbyexamples.com/spark/sp

## 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:/

/www.f (https:/

(https:/   acebo   /www.li (https:/

/twitter   ok.co   nkedin   /github

.com/s   m/spar   .com/i   .com/s

parkby   kbyexa   n/n-nk-   park-

examp   mples/   b860a   examp

les)        )        8193/)   les/)

ark-check-column-data-type-is-integer-
or-string/)

ark-check-column-data-type-is-integer-
or-string/)