

Medium

 Search Write

★ Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#)



# Apache Spark: Explode Function



BigDataEnthusiast

Follow

6 min read · Aug 15, 2023



Apache Spark built-in function that takes input as an column object (array or map type) and returns a new row for each element in the given array or map type column. Refer official documentation [here](#).

```
1 package com
2
3 import org.apache.spark.sql.SparkSession
4 import org.apache.spark.sql.functions.explode
5
6 explode(e: Column) Column
7 explode_outer(e: Column) Column
8 posexplode(e: Column) Column
9 posexplode_outer(e: Column) Column
```

## Explode

Here is the example of `explode` function.

```
1 package com
2
3 import org.apache.spark.sql.SparkSession
```

```
4  import org.apache.spark.sql.functions._
5
6  object InvokeCode extends App {
7
8      val spark: SparkSession = SparkSession
9          .builder()
10         .master("local[1]")
11         .appName("learn")
12         .getOrCreate()
13
14     val sc = spark.sparkContext
15     sc.setLogLevel("OFF")
16
17     import spark.implicits._
18
19     val df = List(
20         ("1", "A", List("1654123", "2654122", "3654121")),
21         ("2", "B", List("1254123", "2354122", "3454121")),
22         ("3", "C", List()),
23         ("4", "D", null)
24     ).toDF("id", "name", "phone_details")
25
26     df.show(false)
27     df.printSchema()
28
29     val df_exploded =
30         df.withColumn("phone_details_exploded", explode($"phone_details"))
31
32     df_exploded.show(false)
33     df_exploded.printSchema()
34 }
```

Input Dataframe: df

```
+---+---+-----+
|id |name|phone_details          |
+---+---+-----+
|1  |A   |[1654123, 2654122, 3654121]|
|2  |B   |[1254123, 2354122, 3454121]|
|3  |C   |[]
```

```
|4 |D |null|
+---+---+-----+-----+

```

Output Dataframe: df\_exploded

```
+---+---+-----+-----+
|id |name|phone_details|phone_details_exploded|
+---+---+-----+-----+
|1 |A |[1654123, 2654122, 3654121]|1654123|
|1 |A |[1654123, 2654122, 3654121]|2654122|
|1 |A |[1654123, 2654122, 3654121]|3654121|
|2 |B |[1254123, 2354122, 3454121]|1254123|
|2 |B |[1254123, 2354122, 3454121]|2354122|
|2 |B |[1254123, 2354122, 3454121]|3454121|
+---+---+-----+-----+

```

If you have noticed here, exploded doesn't include empty list/null values. To include these null values we have to use `explode_outer` function.

```
val df_exploded =
  df.withColumn("phone_details_exploded", explode_outer($"phone_details"))

df_exploded.show(false)
```

Output Dataframe: df\_exploded using explode\_outer

```
+---+---+-----+-----+
|id |name|phone_details|phone_details_exploded|
+---+---+-----+-----+
|1 |A |[1654123, 2654122, 3654121]|1654123|
|1 |A |[1654123, 2654122, 3654121]|2654122|
|1 |A |[1654123, 2654122, 3654121]|3654121|
|2 |B |[1254123, 2354122, 3454121]|1254123|
|2 |B |[1254123, 2354122, 3454121]|2354122|
|2 |B |[1254123, 2354122, 3454121]|3454121|
|3 |C |[ ]|null|
|4 |D |null|null|

```

+---+---+-----+-----+-----+

Let's explode if we have array of objects instead array of strings/integers.

## Explode — Array of Struct

In this example we will explore array of struct.

```
1  package com
2
3  import org.apache.spark.sql.{Row, SparkSession}
4  import org.apache.spark.sql.functions.explode
5  import org.apache.spark.sql.types._
6
7  object InvokeCode extends App {
8
9      val spark: SparkSession = SparkSession
10         .builder()
11         .master("local[1]")
12         .appName("learn")
13         .getOrCreate()
14
15     val sc = spark.sparkContext
16     sc.setLogLevel("OFF")
17
18     import spark.implicits._
19
20     val dataSchema = new StructType()
21         .add("name", StringType)
22         .add("id", StringType)
23         .add(
24             "phone_details",
25             ArrayType(
26                 new StructType()
```

---

```

|name|id |phone_details
+---+---+-----+
|1   |A   |[[home, +1 1254, true], [office, +1 12345, false]]|
|2   |B   |[[home, +1 1254, false], [office, +1 12345, true]]|
|3   |C   |[[home, +1 1254, true], [office, +1 12345, false]]|
|4   |D   |[]
|5   |D   |null
+---+---+-----+

```

Output Dataframe: df\_exploded

```

+---+---+-----+-----+
|name|id |phone_details          |phone_details_e
+---+---+-----+-----+
|1   |A   |[[home, +1 1254, true], [office, +1 12345, false]]|[home, +1 1254,
|1   |A   |[[home, +1 1254, true], [office, +1 12345, false]]|[office, +1 123
|2   |B   |[[home, +1 1254, false], [office, +1 12345, true]]|[home, +1 1254,
|2   |B   |[[home, +1 1254, false], [office, +1 12345, true]]|[office, +1 123
|3   |C   |[[home, +1 1254, true], [office, +1 12345, false]]|[home, +1 1254,
|3   |C   |[[home, +1 1254, true], [office, +1 12345, false]]|[office, +1 123
+---+---+-----+-----+

```

Input Dataframe Schema:

```

root
|-- name: string (nullable = true)
|-- id: string (nullable = true)
|-- phone_details: array (nullable = true)
|   |-- element: struct (containsNull = true)
|   |   |-- phone_type: string (nullable = true)
|   |   |-- number: string (nullable = true)
|   |   |-- primary: boolean (nullable = true)

```

Output Dataframe Schema:

```

root
|-- name: string (nullable = true)
|-- id: string (nullable = true)
|-- phone_details: array (nullable = true)
|   |-- element: struct (containsNull = true)
|   |   |-- phone_type: string (nullable = true)
|   |   |-- number: string (nullable = true)
|   |   |-- primary: boolean (nullable = true)
|-- phone_details_exploded: struct (nullable = true)
|   |-- phone_type: string (nullable = true)
|   |-- number: string (nullable = true)
|   |-- primary: boolean (nullable = true)

```

Notice in the above schema comparison, how the array of struct converted to struct after explode.

Elegant way to select all columns from `struct` see below example.

```
val df_exploded =  
  df.withColumn("phone_details_exploded", explode($"phone_details"))  
  
df_exploded  
  .select("id", "name", "phone_details_exploded.*")  
  .show(false)
```

```
+---+-----+-----+-----+-----+  
|id |name|phone_type|number  |primary|  
+---+-----+-----+-----+-----+  
|A  |1  |home      |+1 1254 |true   |  
|A  |1  |office     |+1 12345|false  |  
|B  |2  |home      |+1 1254 |false  |  
|B  |2  |office     |+1 12345|true   |  
|C  |3  |home      |+1 1254 |true   |  
|C  |3  |office     |+1 12345|false  |  
+---+-----+-----+-----+-----+
```

Other way of selecting columns by explicitly giving column names from struct.

```
df_exploded
  .select(
    "id",
    "phone_details_exploded.phone_type",
    "phone_details_exploded.number"
  )
  .show(false)
```

## Explode Multiple Columns

Suppose we want to explode multiple columns: If we go with one by one approach for exploding multiple columns, it can create bunch of redundant data. We can do first zip columns & then explode. See below example.

```
1 package com
2
3 import org.apache.spark.sql.{Row, SparkSession}
4 import org.apache.spark.sql.functions.{explode, arrays_zip}
5 import org.apache.spark.sql.types._
6
7 object InvokeCode extends App {
8
9     val spark: SparkSession = SparkSession
10       .builder()
11       .master("local[1]")
12       .appName("learn")
13       .getOrCreate()
14
15     val sc = spark.sparkContext
16     sc.setLogLevel("OFF")
17
18     import spark.implicits._
19
20     val dataSchema = new StructType()
21       .add("name", StringType)
22       .add("id", StringType)
```



```

23      .add(
24          "phone_details",
25          ArrayType(
26              new StructType()
27                  .add("phone_type", StringType)
28                  .add("number", StringType)
29                  .add("primary", BooleanType)
30          )
31      )
32      .add(
33          "address_details",
34          ArrayType(
35              new StructType()
36                  .add("address_type", StringType)
37                  .add("address", StringType)
38                  .add("primary", BooleanType)
39          )
40      )
41
42      val data = List(
43          Row(
44              "1",
45              "A",
46              List(
47                  Row("home", "+1 1254", true),
48                  Row("home", "+1 1233", false),
49                  Row("office", "+1 12345", false)
50              ),
51              List(Row("home", "XYZ Lane", true), Row("office", "YZ CITY", false))
52          ),
53          Row(
54              "2",
55              "B",
56              List(Row("home", "+1 1254", false), Row("office", "+1 12345", true)),
57              List(Row("home", "ADC Lane", true), Row("office", "AYZ CITY", false))
58          ),
59          Row(
60              "3",
61              "C",
62              List(Row("home", "+1 1254", true), Row("office", "+1 12345", false)),
63              List(Row("home", "XSD Lane", true), Row("office", "XS CITY", false))
64          ),
65      )

```

```

64     ),
65     Row("4", "D", List(), null),
66     Row("5", "D", null, List())
67 )
68 val rdd = spark.sparkContext.parallelize(data)
69 val df = spark.createDataFrame(rdd, dataSchema)
70
71 df.show(false)
72 df.printSchema()
73
74 val df_exploded = df
75   .withColumn(
76     "zip_col",
77     explode(arrays_zip($"phone_details", $"address_details"))
78   )
79   .drop("phone_details", "address_details")
80   .select("id", "name", "zip_col.*")
81
82 df_exploded.printSchema()
83
84 df_exploded.show(false)

```

Input Dataframe:

```

+----+----+-----+
|name|id |phone_details
+----+----+-----+
|1   |A  |[[home, +1 1254, true], [home, +1 1233, false], [office, +1 12345,
|2   |B  |[[home, +1 1254, false], [office, +1 12345, true]]
|3   |C  |[[home, +1 1254, true], [office, +1 12345, false]]
|4   |D  |[]
|5   |D  |null
+----+----+-----+

```

root

```

|-- name: string (nullable = true)
|-- id: string (nullable = true)
|-- phone_details: array (nullable = true)
|   |-- element: struct (containsNull = true)
|   |   |-- phone_type: string (nullable = true)
|   |   |-- number: string (nullable = true)
|   |   |-- primary: boolean (nullable = true)
|-- address_details: array (nullable = true)

```

```
|      |-- element: struct (containsNull = true)
|      |      |-- address_type: string (nullable = true)
|      |      |-- address: string (nullable = true)
|      |      |-- primary: boolean (nullable = true)
```

Output Dataframe:

```
+---+---+-----+-----+
|id|name|phone_details          |address_details          |
+---+---+-----+-----+
|A|1|[home, +1 1254, true] |[home, XYZ Lane, true] |
|A|1|[home, +1 1233, false] |[office, YZ CITY, false] |
|A|1|[office, +1 12345, false]|null                    |
|B|2|[home, +1 1254, false] |[home, ADC Lane, true] |
|B|2|[office, +1 12345, true] |[office, AYZ CITY, false] |
|C|3|[home, +1 1254, true] |[home, XSD Lane, true] |
|C|3|[office, +1 12345, false] |[office, XS CITY, false] |
+---+---+-----+-----+
```

```
root
|-- id: string (nullable = true)
|-- name: string (nullable = true)
|-- phone_details: struct (nullable = true)
|   |-- phone_type: string (nullable = true)
|   |-- number: string (nullable = true)
|   |-- primary: boolean (nullable = true)
|-- address_details: struct (nullable = true)
|   |-- address_type: string (nullable = true)
|   |-- address: string (nullable = true)
|   |-- primary: boolean (nullable = true)
```

## posexplode

Creates a new row for each element with position in the given array or map column. Uses the default column name `pos` for position, and `col` for elements in the array and `key` and `value` for elements in the map unless specified otherwise.

```
1 package com
2
3 import org.apache.spark.sql.SparkSession
```

```
4  import org.apache.spark.sql.functions._
5
6  object InvokeCode extends App {
7
8      val spark: SparkSession = SparkSession
9          .builder()
10         .master("local[1]")
11         .appName("learn")
12         .getOrCreate()
13
14     val sc = spark.sparkContext
15     sc.setLogLevel("OFF")
16
17     import spark.implicits._
18
19     val df = List(
20         ("1", "A", List("1654123", "2654122", "3654121")),
21         ("2", "B", List("1254123", "2354122", "3454121")),
22         ("3", "C", List()),
23         ("4", "D", null)
24     ).toDF("id", "name", "phone_details")
25
26     df.show(false)
27     df.select($"id", $"name", posexplode($"phone_details")).show(false)
28 }
```

Input Dataframe:

```
+---+-----+-----+
|id |name|phone_details          |
+---+-----+-----+
|1  |A   |[1654123, 2654122, 3654121]|
|2  |B   |[1254123, 2354122, 3454121]|
|3  |C   |[]                          |
|4  |D   |null                        |
+---+-----+-----+
```

Output Dataframe:

```
+---+-----+---+-----+
|id |name|pos|col   |
+---+-----+---+-----+
```

```
|1|A|0|1654123|
|1|A|1|2654122|
|1|A|2|3654121|
|2|B|0|1254123|
|2|B|1|2354122|
|2|B|2|3454121|
+---+-----+---+-----+
```

If you are using `posexplode` in `withColumn` it might fail with this exception.  
Checkout this issue for more details: [SPARK-20174](#)

```
val df_exploded =
  df.withColumn("phone", posexplode($"phone_details"))
```

Exception in thread "main" org.apache.spark.sql.AnalysisException: The number

So better to use `posexplode` with `select` or `selectExpr`.

```
import spark.implicits._

val df = List(
  ("1", "A", List("1654123", "2654122", "3654121")),
  ("2", "B", List("1254123", "2354122", "3454121")),
  ("3", "C", List()),
  ("4", "D", null)
).toDF("id", "name", "phone_details")

df.selectExpr("*, posexplode(phone_details) as (p,c)")
  .drop("phone_details")
  .show(false)
```

**Output:**

```
+---+---+---+-----+
|id |name|p  |c      |
+---+---+---+-----+
|1  |A   |0  |1654123|
|1  |A   |1  |2654122|
|1  |A   |2  |3654121|
|2  |B   |0  |1254123|
|2  |B   |1  |2354122|
|2  |B   |2  |3454121|
+---+---+---+-----+
```

## Refer other blogs:

### Apache Spark — createDataFrame

createDataFrame is an overloaded method present in SparkSession class (org.apache.spark.sql) and there are...

[bigdataenthusiast.medium.com](https://bigdataenthusiast.medium.com)

### Apache Spark — implicits Object (Implicits Conversions)

In the SparkSession class there is one object defined as implicits, which extends SQLImplicits abstract class. So once...

[bigdataenthusiast.medium.com](https://bigdataenthusiast.medium.com)

## References:

- <https://spark.apache.org>
- [SPARK-20174](#)

Apache Spark

Spark Scala Tutorial



## Written by BigDataEnthusiast

91 followers · 4 following

Follow

AWS Certified Data Engineer | Databricks Certified Apache Spark 3.0 Developer |  
Databricks Certified Data Engineer | Oracle Certified SQL Expert

## No responses yet

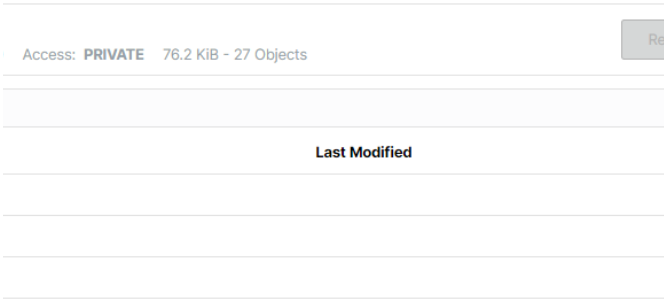



Aaditya Adhikari

What are your thoughts?

## More from BigDataEnthusiast



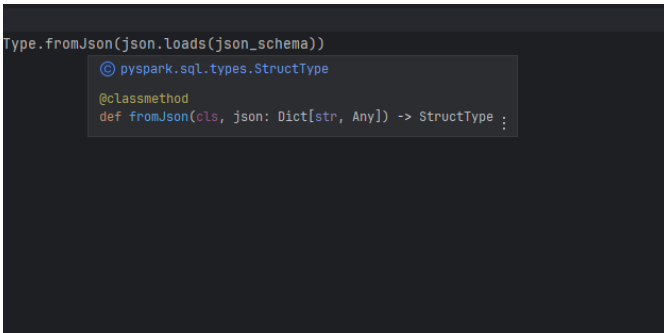


 BigDataEnthusiast

## Apache Iceberg—Hidden Partitioning

In this blog we will explore “Hidden Partitioning” concept in Apache Iceberg.


Mar 31, 2024  74  



 BigDataEnthusiast

## How to create StructType schema from JSON schema | PySpark

Using Apache Spark class `pyspark.sql.types.StructType` method...

Nov 25, 2024  3  

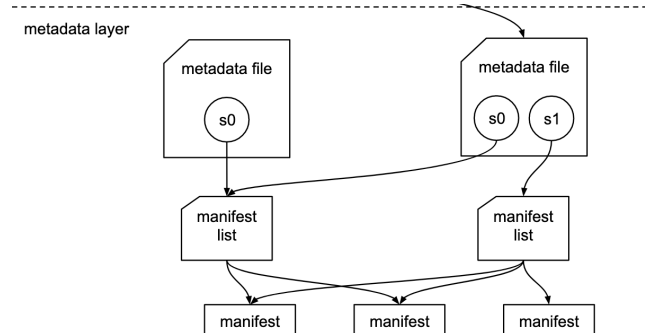


 BigDataEnthusiast

## Apache Iceberg—Insert Overwrite

INSERT OVERWRITE can replace/overwrite the data in iceberg table, depending on...

Jul 9, 2023  5  



 BigDataEnthusiast

## Apache Iceberg Table Format Versions

In this blog we will explore mainly these things.

Jun 25, 2023  10  

See all from BigDataEnthusiast

## Recommended from Medium



 In Art of Data Engineering by Santosh Joshi

 Sriw World of Coding

### Understanding Collect, Take, Limit, Show, Head and Display i...

A Quick and Crisp Guide to Inspecting DataFrames Efficiently in PySpark

### Optimizing Spark Resource Allocation with spark-submit

While there are several ways to optimize Spark jobs, one of the most fundamental...

★ Jan 13 🖱 7

🔖+ ...

May 4 🖱 2

🔖+ ...





Siddharth Ghosh


## Spark Interview Series—Different Optimization Techniques(RBO,...

Apache Spark has evolved into a powerful big data processing engine with several...

 Apr 8





Diogo Santos

## Mastering Apache Spark Shuffle Optimization: Techniques for...

Explore the power of bucketing, repartitioning, and broadcast joins to...

 Jan 26  64





Krishnatej Sreeramula

## SQL Made Simple: Spark SQL Basics 🚀💻

Introduction

 Jan 5  3


Overwrite Entire Table	Small tables	Simple	Inefficient for large tables
Update Specific Partitions	Partitioned tables	Efficient for large tables	Limited to partition-level updates
Merge Using Delta Lake	Frequent updates, ACID compliance	Scalable, ACID-compliant	Requires Delta Lake
Use a Join to Update Records	Standard Spark tables	Flexible	Requires overwriting the entire table
Use a Temporary Table	Atomicity without Delta Lake	Ensures atomicity	Requires additional storage
SQL UPDATE	Delta Lake, SQL-based	Simple and	Requires Delta Lake




Sujatha Mudadla

## Section 2: Data Processing—Updating Records in a Spark...

In data processing, Type 1 updates refer to overwriting existing records with new data...

 Jan 11  1

See more recommendations

---

[Help](#) [Status](#) [About](#) [Careers](#) [Press](#) [Blog](#) [Privacy](#) [Rules](#) [Terms](#) [Text to speech](#)