

1. [Apache Spark Professional Training with Hands On Lab Sessions](#)
2. [Oreilly Databricks Apache Spark Developer Certification Simulator](#)

SPARK ADVANCED : DATA PARTITIONING

By www.HadoopExam.com

Note: These instructions should be used with the HadoopExam Apache Spark: Professional Trainings.
Where it is executed and you can do hands on with trainer.

1. Hadoop Training
2. Spark Training
3. HBase Training
4. MapR Developer
5. MapR HBase
6. CCA500 Certification
7. Spark Certification
8. EMC Data Science

Hadoop Specialization offer == 50% + 35% off

Hadoop Expert

~~52000INR~~ == 16900INR Only
~~\$1150~~ == \$373 Only
[Hadoop Specialization offer](#)

* @ End of the Offer Prices will increase by 25%

Limited Time Offer (Less Than 5Days Remain)



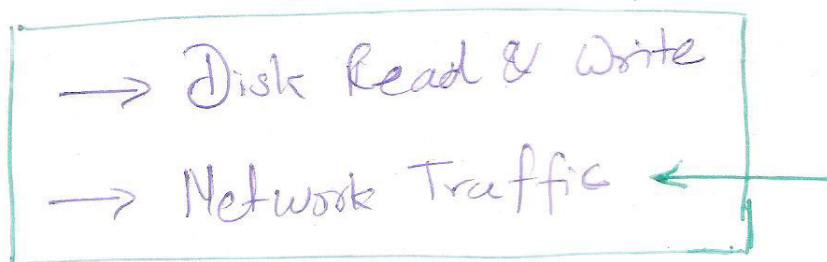
[Cloudera CCA175 \(Hadoop and Spark Developer Hands-on Certification available with total 75 solved problem scenarios. Click for More Detail\)](#)

1. What is Partitioning and why?
2. Data Partitioning example using Join (Hash Partitioning)
3. Understand Partitioning using Example for get Recommendations for Customer
4. Understand Partitioning code using Spark-Scala
5. Operations which create Partitioned RDD
6. Operation which get benefit of Partitioning
7. Operation that affect the partitioning

Spark: Data Partitioning

⇒ Distributing Data Across nodes in cluster, to improve efficiency.

⇒ To biggest performance constraint, while working on distributed platform.



⇒ Data partitioning will reduce n/w traffic.

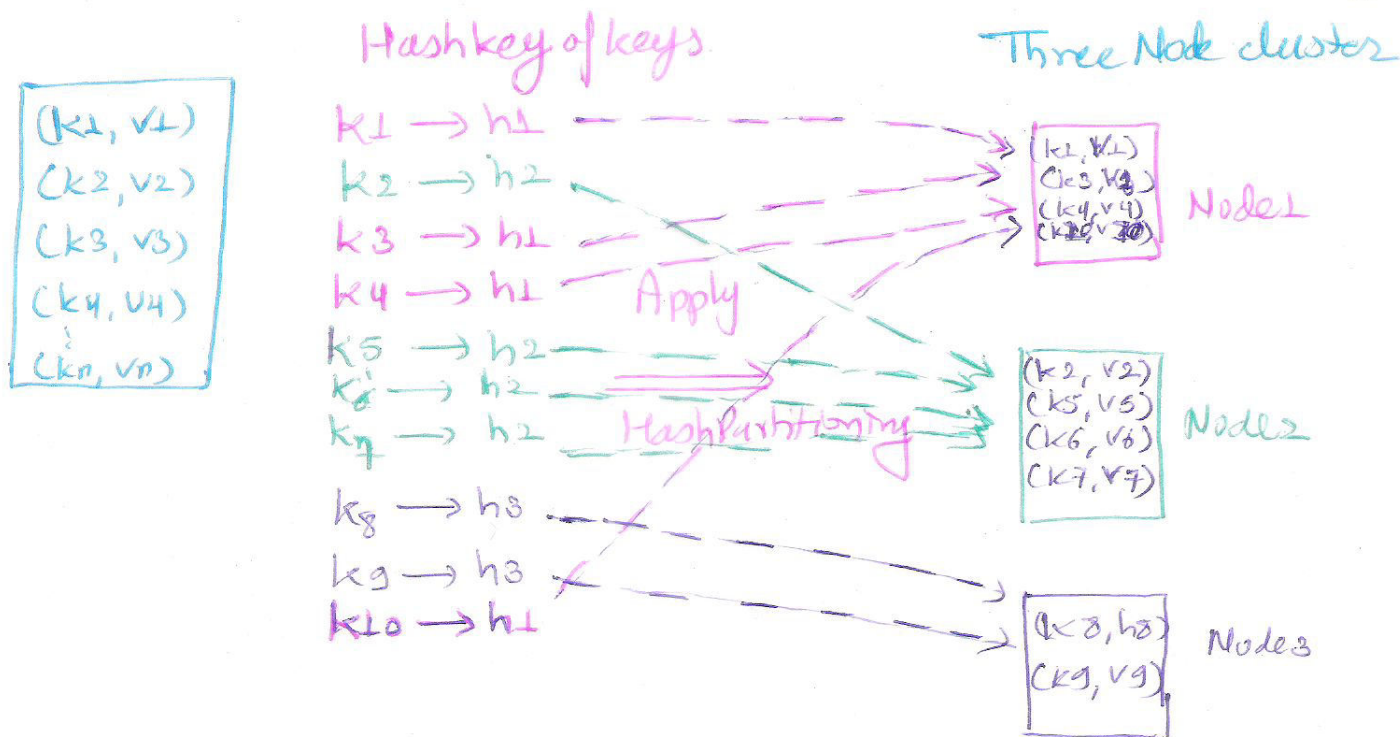
⇒ Always useful: No

⇒ If given RDD is scanned only once, then no use of partitioning

⇒ When ⇒ In key-oriented operations (maybe PairRDD), and dataset is re-used multiple times. e.g. Lookup Data, used in Join operation.

⇒ Hashkey based partitioning: -

(2)



⇒ Now calling persist() method on RDD is Mandatory

Example:- Let's take an example from ~~an~~ Amazon kindle online, web reading application.

~~Customer~~

Customer-Info (RDD structure)

```
[CustomerId, List [book1, book2, book3 ---- bookn]
  1, [book1, book2, book7, book100 ---- book700]
  2, [book2, book9, book400 ---- book405]
  ----
  1001, [book7, book8, book4 ---- book1001]
```

⇒ CustomerInfo RDD is huge in size.

③

⇒ Another RDD, which we receive every 10 minutes using Streaming Application. With following structure.

CustomerId, bookVisited	
1,	book1
7,	book1001
8,	book44
9,	book1001
20,	book11001
⋮	
1001,	book72

www.HadoopExam.com

Problem Statement ⇒ Customer-Interest

Find out the book name, which is not subscribed by a Customer, but he had visited.

For Example: - CustomerId → 1001
Book Visited → book72 (Not subscribed)

Join both the RDD, so Amazon® can find, which book is not subscribed by customer, but he has interest so it can send recommendation in his kindle® box.

⇒ Why partitioning will help here.

Partitioning of Customer-Info RDD, is required. Because this RDD, we will be using in every 10 minutes to calculate customer Interest.

Two RDDs

Customer-Info

Big Table
(No need to re-calculate again and again)

Customer-Interest

Re-generated every 10 minutes from Streaming application.

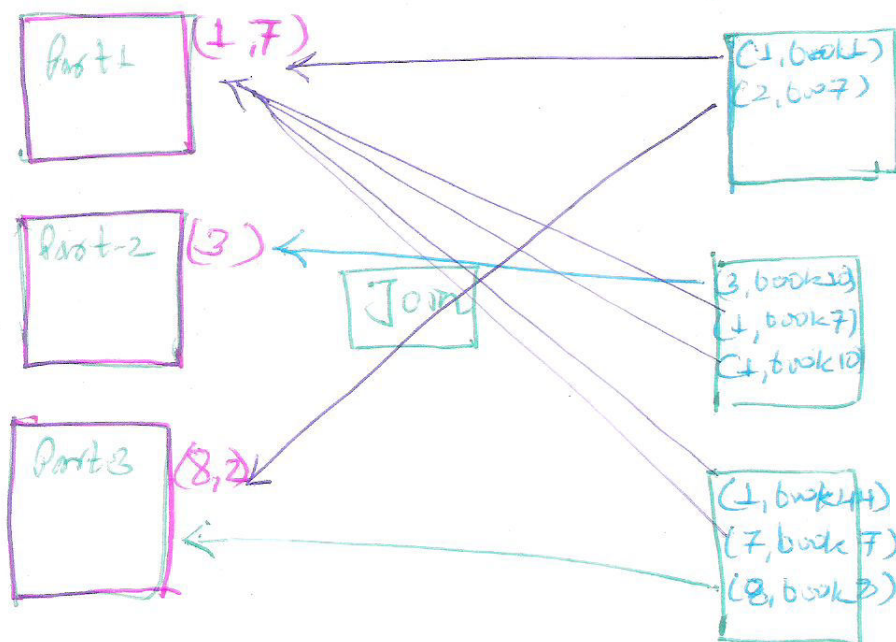


Partition this RDD

& persist() its state

Customer-Info (RDD)

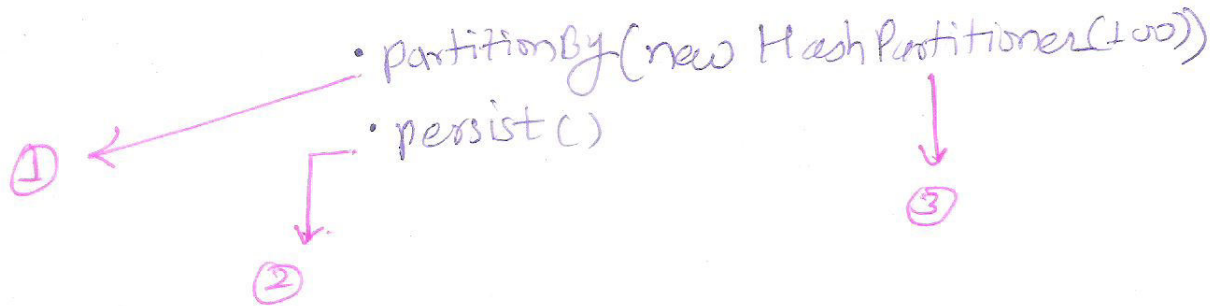
Customer-Interest (RDD)



Application Code Looklike

~~val user~~

val customerData = sc.textFile [custId, list [books]]



- 1 → partitionBy is a method to partition data.
- 2 → Algorithm (technique) used to do partitioning
- 3 → Save the partition created (Don't miss it)

```
def generateData (streamedfile: String) {
```

```
  val row = sc.textFile [custId, book] (streamfile)
```

```
  val result = customerData.join(row)
```

```
  // (CustId, [List(Book), Book]) RDD structure joined
```

```
  val newVisits = result.filter {
```

```
    Case (custId, (list, book)) =>
```

```
      ! list.books.contains (book.name)
```

```
  }
```

```
  newVisits.saveAsTextFile ("path"),
```

⇒ Operations create partitioned RDD :-

There are many operations which will create partitioned RDD.

sortByKey() → Range Partition.

groupByKey() → hash-partitioned RDD.

map() → Forget Parent RDD partitions

Why? ⇒ This map() operation can modify the key

⇒ Operations that Benefit from Partitioning: -

→ Operations which requires shuffling data by key across network.

→

coGroup()	rightOuterJoin()	lookup()
groupWith()	groupByKey()	
join()	reduceByKey()	
leftOuterJoin()	combineByKey()	

→ reduceByKey() :- Will work only on single RDD.

k1
k7
k8
k9

Part-1 (do operation locally in
① each partition)

k2
k3
k4
k5
k6

Part-2 ② Send Result to master node
from each node.

Operation That affect Partitioning:-

Calling `reduceByKey()` after `join()`

Significantly faster.

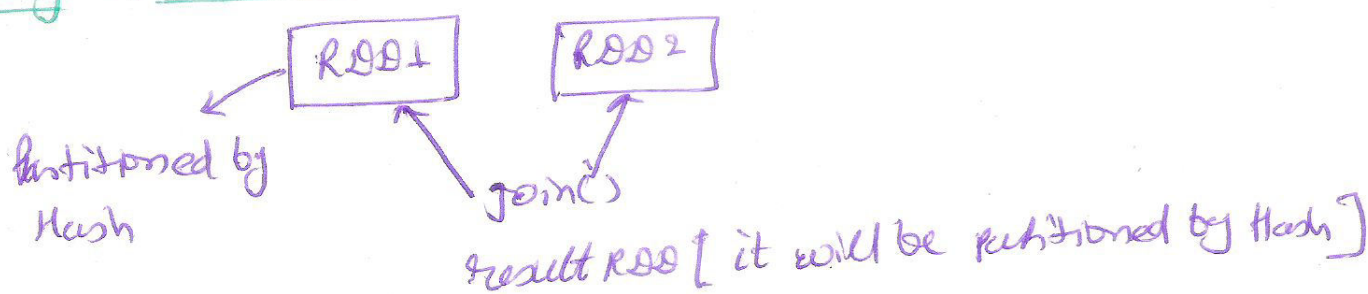
Why? `join` will create hash-partitioned RDD and `reduceByKey()` will take benefit of this.

=> `map()` can change key.

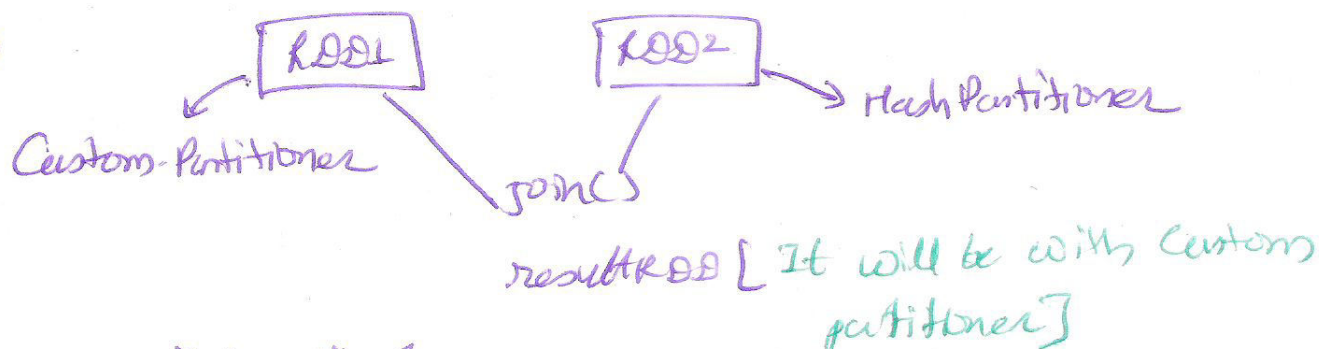
=> `mapValues()` & `flatMapValues()`: will not change the key.

Binary Operation:-

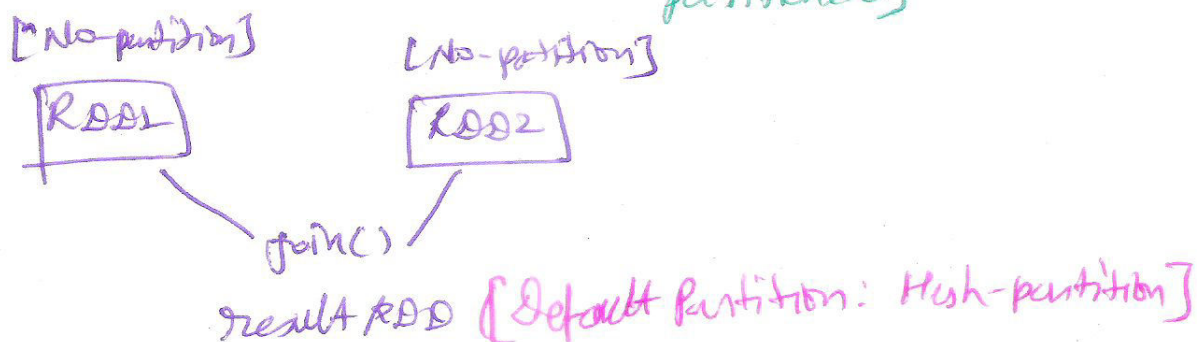
①



②



③



cloudera CCA175 Is Now Available , with Hands On Sessions

CCA Spark and Hadoop Developer Certification

HadoopExam Learning Resource provides the following material for the Advanced Technologies.
Please visit www.HadoopExam.com for more detail this is just a few products from portfolio.

Price start for training with Just \$79/3500INR



Apache Spark
Professional
Training with
HandsOn Session

+ Certification
Material



Hadoop Professional
Training with
HandsOn Session

+ Certification
Material



HBase Professional
Training with
HandsOn Session

+ Certification
Material



Certification
Material



Certification
Material



Certification
Material



Certification
Material



Certification
Material



Microsoft Azure

Certification
Material



Certification
Material