

# Complete PySpark RDD Methods Reference Guide

## TRANSFORMATIONS (Lazy - Return new RDD)

### Basic Transformations

Method	Returns	Description	Example
<code>map(func)</code>	RDD	Apply function to each element	<code>rdd.map(lambda x: x*2)</code>
<code>filter(func)</code>	RDD	Filter elements by condition	<code>rdd.filter(lambda x: x &gt; 5)</code>
<code>flatMap(func)</code>	RDD	Apply function and flatten result	<code>rdd.flatMap(lambda x: x.split())</code>
<code>distinct()</code>	RDD	Remove duplicate elements	<code>rdd.distinct()</code>
<code>sample(withReplacement, fraction, seed)</code>	RDD	Random sample of elements	<code>rdd.sample(False, 0.1)</code>

### Set Operations

Method	Returns	Description	Example
<code>union(other)</code>	RDD	Union with another RDD	<code>rdd1.union(rdd2)</code>
<code>intersection(other)</code>	RDD	Common elements between RDDs	<code>rdd1.intersection(rdd2)</code>
<code>subtract(other)</code>	RDD	Elements in this RDD but not other	<code>rdd1.subtract(rdd2)</code>
<code>cartesian(other)</code>	RDD	Cartesian product with another RDD	<code>rdd1.cartesian(rdd2)</code>

### Key-Value Transformations (for Pair RDDs)

Method	Returns	Description	Example
<code>mapValues(func)</code>	RDD	Apply function only to values	<code>pair_rdd.mapValues(lambda x: x*2)</code>
<code>keys()</code>	RDD	Extract keys only	<code>pair_rdd.keys()</code>
<code>values()</code>	RDD	Extract values only	<code>pair_rdd.values()</code>
<code>reduceByKey(func)</code>	RDD	Reduce values by key	<code>pair_rdd.reduceByKey(lambda a,b: a+b)</code>

Method	Returns	Description	Example
<code>groupByKey()</code>	RDD	Group values by key	<code>pair_rdd.groupByKey()</code>
<code>sortByKey(ascending)</code>	RDD	Sort by key	<code>pair_rdd.sortByKey()</code>
<code>join(other)</code>	RDD	Inner join on keys	<code>rdd1.join(rdd2)</code>
<code>leftOuterJoin(other)</code>	RDD	Left outer join	<code>rdd1.leftOuterJoin(rdd2)</code>
<code>rightOuterJoin(other)</code>	RDD	Right outer join	<code>rdd1.rightOuterJoin(rdd2)</code>
<code>fullOuterJoin(other)</code>	RDD	Full outer join	<code>rdd1.fullOuterJoin(rdd2)</code>
<code>cogroup(other)</code>	RDD	Group together values from both RDDs	<code>rdd1.cogroup(rdd2)</code>
<code>subtractByKey(other)</code>	RDD	Remove elements with keys in other	<code>rdd1.subtractByKey(rdd2)</code>

## Partitioning & Ordering

Method	Returns	Description	Example
<code>partitionBy(numPartitions, partitioner)</code>	RDD	Custom partitioning	<code>rdd.partitionBy(4)</code>
<code>repartition(numPartitions)</code>	RDD	Repartition with shuffle	<code>rdd.repartition(8)</code>
<code>coalesce(numPartitions)</code>	RDD	Reduce partitions without shuffle	<code>rdd.coalesce(2)</code>
<code>sortBy(func, ascending)</code>	RDD	Sort by function result	<code>rdd.sortBy(lambda x: x[1])</code>
<code>randomSplit(weights, seed)</code>	List[RDD]	Split RDD randomly	<code>rdd.randomSplit([0.7, 0.3])</code>

## Advanced Transformations

Method	Returns	Description	Example
<code>mapPartitions(func)</code>	RDD	Apply function to each partition	<code>rdd.mapPartitions(process_partition)</code>
<code>mapPartitionsWithIndex(func)</code>	RDD	Map partitions with index	<code>rdd.mapPartitionsWithIndex(func)</code>
<code>glom()</code>	RDD	Convert each partition to	<code>rdd.glom()</code>

Method	Returns	Description	Example
<code>pipe(command)</code>	RDD	Pipe through external command	<code>rdd.pipe("grep pattern")</code>
<code>keyBy(func)</code>	RDD	Create key-value pairs	<code>rdd.keyBy(lambda x: x[0])</code>
<code>zipWithIndex()</code>	RDD	Zip with element indices	<code>rdd.zipWithIndex()</code>
<code>zipWithUniquelId()</code>	RDD	Zip with unique IDs	<code>rdd.zipWithUniquelId()</code>
<code>zip(other)</code>	RDD	Zip with another RDD	<code>rdd1.zip(rdd2)</code>

## ⚡ ACTIONS (Eager - Trigger computation, return values to driver)

### Collection Actions

Method	Returns	Description	Example
<code>collect()</code>	List	Return all elements as list	<code>rdd.collect()</code>
<code>take(n)</code>	List	Return first n elements	<code>rdd.take(10)</code>
<code>takeOrdered(n, key)</code>	List	Return n smallest elements	<code>rdd.takeOrdered(5, key=lambda x: x)</code>
<code>top(n, key)</code>	List	Return n largest elements	<code>rdd.top(5)</code>
<code>takeSample(withReplacement, n, seed)</code>	List	Return random sample	<code>rdd.takeSample(False, 10)</code>
<code>first()</code>	Element	Return first element	<code>rdd.first()</code>

## Aggregation Actions

Method	Returns	Description	Example
<code>reduce(func)</code>	Element	Reduce elements using function	<code>rdd.reduce(lambda a,b: a+b)</code>
<code>fold(zeroValue, func)</code>	Element	Fold with initial value	<code>rdd.fold(0, lambda a,b: a+b)</code>
<code>aggregate(zeroValue, seqOp, combOp)</code>	Element	General aggregation	<code>rdd.aggregate(0, add, add)</code>
<code>sum()</code>	Number	Sum of elements	<code>rdd.sum()</code>
<code>count()</code>	Int	Number of elements	<code>rdd.count()</code>
<code>countByValue()</code>	Dict	Count occurrences of each value	<code>rdd.countByValue()</code>
<code>countByKey()</code>	Dict	Count occurrences by key	<code>pair_rdd.countByKey()</code>
<code>max()</code>	Element	Maximum element	<code>rdd.max()</code>
<code>min()</code>	Element	Minimum element	<code>rdd.min()</code>
<code>mean()</code>	Float	Average of elements	<code>rdd.mean()</code>
<code>variance()</code>	Float	Variance of elements	<code>rdd.variance()</code>
<code>stdev()</code>	Float	Standard deviation	<code>rdd.stdev()</code>
<code>stats()</code>	StatCounter	Statistical summary	<code>rdd.stats()</code>
<code>histogram(buckets)</code>	Tuple	Histogram of values	<code>rdd.histogram(10)</code>

## Key-Value Actions (for Pair RDDs)

Method	Returns	Description	Example
<code>collectAsMap()</code>	Dict	Collect as dictionary	<code>pair_rdd.collectAsMap()</code>
<code>lookup(key)</code>	List	Values for specific key	<code>pair_rdd.lookup("key1")</code>
<code>countByKey()</code>	Dict	Count by key	<code>pair_rdd.countByKey()</code>

## Output Actions

Method	Returns	Description	Example
<code>saveAsTextFile(path)</code>	None	Save as text files	<code>rdd.saveAsTextFile("output/")</code>

Method	Returns	Description	Example
<code>saveAsPickleFile(path)</code>	None	Save as pickle files	<code>rdd.saveAsPickleFile("output.pkl")</code>
<code>saveAsSequenceFile(path)</code>	None	Save as sequence files	<code>rdd.saveAsSequenceFile("output/")</code>
<code>foreach(func)</code>	None	Apply function to each element	<code>rdd.foreach(print)</code>
<code>foreachPartition(func)</code>	None	Apply function to each partition	<code>rdd.foreachPartition(process)</code>

## Boolean Actions

Method	Returns	Description	Example
<code>isEmpty()</code>	Boolean	Check if RDD is empty	<code>rdd.isEmpty()</code>



## UTILITY METHODS (Metadata & Control)

### RDD Information

Method	Returns	Description	Example
<code>getNumPartitions()</code>	Int	Number of partitions	<code>rdd.getNumPartitions()</code>
<code>partitions()</code>	List	List of partitions	<code>rdd.partitions()</code>
<code>partitioner()</code>	Partitioner	Partitioner object	<code>rdd.partitioner()</code>
<code>glom()</code>	RDD[List]	Show partition contents	<code>rdd.glom().collect()</code>
<code>toDebugString()</code>	String	Debug information	<code>rdd.toDebugString()</code>

### Caching & Persistence

Method	Returns	Description	Example
<code>cache()</code>	RDD	Cache in memory	<code>rdd.cache()</code>
<code>persist(storageLevel)</code>	RDD	Persist with storage level	<code>rdd.persist(MEMORY_AND_DISK)</code>
<code>unpersist(blocking)</code>	RDD	Remove from cache	<code>rdd.unpersist()</code>
<code>checkpoint()</code>	None	Checkpoint RDD	<code>rdd.checkpoint()</code>
<code>localCheckpoint()</code>	RDD	Local checkpoint	<code>rdd.localCheckpoint()</code>
<code>isCheckpointed()</code>	Boolean	Check if checkpointed	<code>rdd.isCheckpointed()</code>

Method	Returns	Description	Example
<code>getCheckpointFile()</code>	String	Get checkpoint file	<code>rdd.getCheckpointFile()</code>

## Dependencies & Lineage

Method	Returns	Description	Example
<code>dependencies()</code>	List	RDD dependencies	<code>rdd.dependencies()</code>
<code>context()</code>	SparkContext	Get SparkContext	<code>rdd.context()</code>
<code>name()</code>	String	RDD name	<code>rdd.name()</code>
<code>setName(name)</code>	RDD	Set RDD name	<code>rdd.setName("my_rdd")</code>
<code>id()</code>	Int	Unique RDD ID	<code>rdd.id()</code>

## STORAGE LEVELS (for persist())

Storage Level	Memory	Disk	Serialized	Replicated
<code>MEMORY_ONLY</code>	✓	✗	✗	✗
<code>MEMORY_AND_DISK</code>	✓	✓	✗	✗
<code>MEMORY_ONLY_SER</code>	✓	✗	✓	✗
<code>MEMORY_AND_DISK_SER</code>	✓	✓	✓	✗
<code>DISK_ONLY</code>	✗	✓	✗	✗
<code>MEMORY_ONLY_2</code>	✓	✗	✗	✓
<code>MEMORY_AND_DISK_2</code>	✓	✓	✗	✓

## METHOD CATEGORIES SUMMARY

### Transformations (Lazy)

- **Element-wise:** `map`, `filter`, `flatMap`, `distinct`
- **Set operations:** `union`, `intersection`, `subtract`
- **Key-Value:** `reduceByKey`, `groupByKey`, `join`, `mapValues`
- **Repartitioning:** `repartition`, `coalesce`, `partitionBy`
- **Sorting:** `sortBy`, `sortByKey`

## ⚡ Actions (Eager)

- **Collection:** `collect`, `take`, `first`, `takeOrdered`
- **Aggregation:** `reduce`, `count`, `sum`, `mean`, `max`, `min`
- **Counting:** `countByValue`, `countByKey`
- **Output:** `saveAsTextFile`, `foreach`
- **Lookup:** `lookup` (for pair RDDs)

## 🔧 Utilities

- **Info:** `count`, `getNumPartitions`, `toDebugString`
  - **Caching:** `cache`, `persist`, `unpersist`
  - **Checkpointing:** `checkpoint`, `isCheckpointed`
- 

## 💡 Performance Tips

1. Use transformations over actions when possible (lazy evaluation)
2. Cache frequently used RDDs with `cache()` or `persist()`
3. Avoid `collect()` on large datasets
4. Use `reduceByKey` over `groupByKey` for aggregations
5. Partition appropriately - too few = underutilization, too many = overhead
6. Use `coalesce` over `repartition` when reducing partitions