*Table 3-2. Basic RDD transformations on an RDD containing {1, 2, 3, 3}*

| Function name | Purpose | Example | Result |
|---|---|---|---|
| map() | Apply a function to each element in the RDD and return an RDD of the result. | rdd.map(x => x + 1) | {2, 3, 4, 4} |
| flatMap() | Apply a function to each element in the RDD and return an RDD of the contents of the iterators returned. Often used to extract words. | rdd.flatMap(x => x.to(3)) | {1, 2, 3, 2, 3, 3, 3} |
| filter() | Return an RDD consisting of only elements that pass the condition passed to filter(). | rdd.filter(x => x != 1) | {2, 3, 3} |
| distinct() | Remove duplicates. | rdd.distinct() | {1, 2, 3} |
| sample(withRe placement, frac tion, [seed]) | Sample an RDD, with or without replacement. | rdd.sample(false, 0.5) | Nondeterministic |

*Table 3-3. Two-RDD transformations on RDDs containing {1, 2, 3} and {3, 4, 5}*

| Function name | Purpose | Example | Result |
|---|---|---|---|
| union() | Produce an RDD containing elements from both RDDs. | rdd.union(other) | {1, 2, 3, 3, 4, 5} |
| intersec tion() | RDD containing only elements found in both RDDs. | rdd.intersection(other) | {3} |
| subtract() | Remove the contents of one RDD (e.g., remove training data). | rdd.subtract(other) | {1, 2} |
| cartesian() | Cartesian product with the other RDD. | rdd.cartesian(other) | {(1, 3), (1, 4), … (3,5)} |

*Table 3-4. Basic actions on an RDD containing {1, 2, 3, 3}*

| Function name | Purpose | Example | Result |
|---|---|---|---|
| `collect()` | Return all elements from the RDD. | `rdd.collect()` | `{1, 2, 3, 3}` |
| `count()` | Number of elements in the RDD. | `rdd.count()` | `4` |
| `countByValue()` | Number of times each element occurs in the RDD. | `rdd.countByValue()` | `{(1, 1), (2, 1), (3, 2)}` |

| Function name | Purpose | Example | Result |
|---|---|---|---|
| `take(num)` | Return num elements from the RDD. | `rdd.take(2)` | `{1, 2}` |
| `top(num)` | Return the top num elements the RDD. | `rdd.top(2)` | `{3, 3}` |
| `takeOrdered(num)(order ing)` | Return num elements based on provided ordering. | `rdd.takeOrdered(2) (myOrdering)` | `{3, 3}` |
| `takeSample(withReplace ment, num, [seed])` | Return num elements at random. | `rdd.takeSample(false, 1)` | Nondeterministic |
| `reduce(func)` | Combine the elements of the RDD together in parallel (e.g., sum). | `rdd.reduce((x, y) => x + y)` | 9 |
| `fold(zero)(func)` | Same as `reduce()` but with the provided zero value. | `rdd.fold(0)((x, y) => x + y)` | 9 |
| `aggregate(zeroValue) (seqOp, combOp)` | Similar to `reduce()` but used to return a different type. | `rdd.aggregate((0, 0)) ((x, y) => (x._1 + y, x._2 + 1), (x, y) => (x._1 + y._1, x._2 + y._2))` | (9, 4) |
| `foreach(func)` | Apply the provided function to each element of the RDD. | `rdd.foreach(func)` | Nothing |