# 06 RDD Operations

# RDD Operations in PySpark

## Introduction to RDDs

Resilient Distributed Datasets (RDDs) form the fundamental building blocks of Apache Spark. Think of an RDD as a distributed collection of data that can be processed in parallel across multiple machines. The "resilient" part means that if any part of your data is lost due to hardware failure, Spark can automatically recover it.

RDDs are immutable, which means once created, they cannot be changed. Instead, operations on RDDs create new RDDs. This immutability is crucial for fault tolerance and enables Spark to track the lineage of transformations, allowing it to rebuild lost partitions.

## Types of RDD Operations

RDD operations fall into two main categories, and understanding this distinction is fundamental to working effectively with Spark:

## 1. Transformations (Lazy Operations)

Transformations are operations that create a new RDD from an existing one. They are "lazy," meaning they don't execute immediately when called. Instead, Spark builds up a computation graph and only executes it when an action is triggered. This lazy evaluation allows Spark to optimize the entire computation pipeline.

## 2. Actions (Eager Operations)

Actions are operations that trigger the execution of transformations and return results to the driver program or write data to external storage. When you call an action, Spark executes all the transformations in the lineage to produce the final result.

## Transformation Operations

Let's explore the most important transformation operations with detailed examples:

## Basic Transformations

```python
from pyspark import SparkContext, SparkConf

# Initialize Spark context
conf = SparkConf().setAppName("RDD_Operations_Demo")
sc = SparkContext(conf=conf)

# Create an RDD from a list
numbers = sc.parallelize([1, 2, 3, 4, 5, 6, 7, 8, 9, 10])
```

## map() - One-to-One Transformation

The map transformation applies a function to each element in the RDD and returns a new RDD with the transformed elements.

```python
# Square each number
squared_numbers = numbers.map(lambda x: x ** 2)
print("Original:", numbers.collect())  # [1, 2, 3, 4, 5, 6, 7, 8, 9, 10]
print("Squared:", squared_numbers.collect())  # [1, 4, 9, 16, 25, 36, 49, 64, 81, 100]


# More complex transformation with string manipulation
words = sc.parallelize(["hello", "world", "spark", "rdd"])
upper_words = words.map(lambda word: word.upper())
print("Uppercase:", upper_words.collect())  # ['HELLO', 'WORLD', 'SPARK', 'RDD']
```

## filter() - Selective Transformation

The filter transformation returns a new RDD containing only elements that satisfy a given condition.

```python
# Filter even numbers
even_numbers = numbers.filter(lambda x: x % 2 == 0)
print("Even numbers:", even_numbers.collect())  # [2, 4, 6, 8, 10]


# Filter words with more than 4 characters
long_words = words.filter(lambda word: len(word) > 4)
print("Long words:", long_words.collect())  # ['hello', 'world', 'spark']
```

## flatMap() - One-to-Many Transformation

The flatMap transformation is similar to map, but each input element can be mapped to zero or more output elements. The results are flattened into a single RDD.

```python
# Split sentences into words
sentences = sc.parallelize(["Hello world", "Apache Spark", "RDD operations"])
words_flat = sentences.flatMap(lambda sentence: sentence.split(" "))
print("Flattened words:", words_flat.collect())
# ['Hello', 'world', 'Apache', 'Spark', 'RDD', 'operations']


# Compare with map (which would create nested structure)
words_nested = sentences.map(lambda sentence: sentence.split(" "))
print("Nested structure:", words_nested.collect())
# [['Hello', 'world'], ['Apache', 'Spark'], ['RDD', 'operations']]
```

## Set Operations

RDDs support mathematical set operations that are particularly useful for data analysis:

```python
# Create two RDDs for set operations
rdd1 = sc.parallelize([1, 2, 3, 4, 5])
rdd2 = sc.parallelize([4, 5, 6, 7, 8])
```

```python
# Union - combines all elements from both RDDs
union_result = rdd1.union(rdd2)
print("Union:", union_result.collect())  # [1, 2, 3, 4, 5, 4, 5, 6, 7, 8]

# Intersection - elements present in both RDDs
intersection_result = rdd1.intersection(rdd2)
print("Intersection:", intersection_result.collect())  # [4, 5]

# Subtract - elements in first RDD but not in second
subtract_result = rdd1.subtract(rdd2)
print("Subtract:", subtract_result.collect())  # [1, 2, 3]

# Remove duplicates
distinct_union = rdd1.union(rdd2).distinct()
print("Distinct Union:", distinct_union.collect())  # [1, 2, 3, 4, 5, 6, 7, 8]
```

## Advanced Transformations

### groupByKey() and reduceByKey()

These operations work with key-value pair RDDs and are fundamental for aggregation operations:

```python
# Create a key-value RDD representing sales by product
sales_data = sc.parallelize([
    ("laptop", 1000), ("mouse", 25), ("laptop", 1200),
    ("keyboard", 75), ("mouse", 30), ("laptop", 800)
])

# groupByKey groups values by key
grouped_sales = sales_data.groupByKey()
# Convert to regular Python data structures for display
grouped_result = grouped_sales.mapValues(list).collect()
print("Grouped sales:", grouped_result)
# [('laptop', [1000, 1200, 800]), ('mouse', [25, 30]), ('keyboard', [75])]

# reduceByKey applies a reduction function to values with the same key
total_sales = sales_data.reduceByKey(lambda x, y: x + y)
print("Total sales by product:", total_sales.collect())
# [('laptop', 3000), ('mouse', 55), ('keyboard', 75)]
```

The key difference between these operations is efficiency: reduceByKey performs local aggregation before shuffling data across the network, making it much faster for large datasets than groupByKey followed by a reduction.

### join() Operations

Join operations combine two key-value RDDs based on their keys:

```python
# Create RDDs for product information
products = sc.parallelize([("laptop", "Electronics"), ("mouse", "Electronics"),
("book", "Media")])
prices = sc.parallelize([("laptop", 1000), ("mouse", 25), ("tablet", 300)])

# Inner join - only keys present in both RDDs
inner_joined = products.join(prices)
print("Inner join:", inner_joined.collect())
# [('laptop', ('Electronics', 1000)), ('mouse', ('Electronics', 25))]

# Left outer join - all keys from left RDD
left_joined = products.leftOuterJoin(prices)
print("Left outer join:", left_joined.collect())
# [('laptop', ('Electronics', 1000)), ('mouse', ('Electronics', 25)), ('book',
('Media', None))]

# Right outer join - all keys from right RDD
right_joined = products.rightOuterJoin(prices)
print("Right outer join:", right_joined.collect())
# [('laptop', ('Electronics', 1000)), ('mouse', ('Electronics', 25)), ('tablet',
(None, 300))]
```

## Action Operations

Actions trigger the execution of transformations and return results. Here are the most commonly used actions:

## Basic Actions

```python
# Create sample data for actions
sample_numbers = sc.parallelize([1, 2, 3, 4, 5, 4, 3, 2, 1])

# collect() - returns all elements as a list (use carefully with large datasets)
all_elements = sample_numbers.collect()
print("All elements:", all_elements)  # [1, 2, 3, 4, 5, 4, 3, 2, 1]

# count() - returns the number of elements
element_count = sample_numbers.count()
print("Count:", element_count)  # 9

# first() - returns the first element
first_element = sample_numbers.first()
print("First element:", first_element)  # 1

# take(n) - returns first n elements
first_three = sample_numbers.take(3)
```

```python
print("First three:", first_three)  # [1, 2, 3]


# top(n) - returns top n elements in descending order
top_three = sample_numbers.top(3)
print("Top three:", top_three)  # [5, 4, 4]
```

## Aggregation Actions

```python
# Mathematical aggregations
sum_result = sample_numbers.sum()
print("Sum:", sum_result)  # 25


# reduce() - applies a function to reduce all elements to a single value
max_value = sample_numbers.reduce(lambda x, y: x if x > y else y)
print("Maximum:", max_value)  # 5


# fold() - similar to reduce but with an initial value
sum_with_initial = sample_numbers.fold(10, lambda x, y: x + y)
print("Sum with initial 10:", sum_with_initial)  # 35


# aggregate() - more flexible aggregation with different types
# Format: aggregate(zeroValue, seqOp, combOp)
stats = sample_numbers.aggregate(
    (0, 0),  # Initial value: (sum, count)
    lambda acc, value: (acc[0] + value, acc[1] + 1),  # Sequential operation
    lambda acc1, acc2: (acc1[0] + acc2[0], acc1[1] + acc2[1])  # Combine operation
)
average = stats[0] / stats[1]
print(f"Average: {average}")  # 2.777...
```

## Sampling and Statistical Actions

```python
# Create larger dataset for sampling
large_numbers = sc.parallelize(range(1, 1001))  # 1 to 1000


# sample() - returns a random sample
# Parameters: withReplacement, fraction, seed
sample_data = large_numbers.sample(False, 0.1, 42)  # 10% sample without
replacement
print("Sample size:", sample_data.count())  # Around 100 elements


# takeSample() - returns exact number of samples
exact_sample = large_numbers.takeSample(False, 10, 42)  # Exactly 10 samples
print("Exact 10 samples:", exact_sample)
```

```python
# Statistical operations
print("Min:", large_numbers.min())  # 1
print("Max:", large_numbers.max())  # 1000
print("Standard deviation:", large_numbers.stdev())  # Standard deviation of 1-
1000
```

## Working with Key-Value Pairs

Key-value operations are essential for many data processing tasks:

```python
# Create word count example - a classic MapReduce problem
text = "the quick brown fox jumps over the lazy dog the fox is quick"
words_rdd = sc.parallelize(text.split())

# Step 1: Map each word to (word, 1)
word_pairs = words_rdd.map(lambda word: (word, 1))
print("Word pairs:", word_pairs.collect())

# Step 2: Reduce by key to count occurrences
word_counts = word_pairs.reduceByKey(lambda x, y: x + y)
print("Word counts:", word_counts.collect())
# [('the', 3), ('quick', 2), ('brown', 1), ('fox', 2), ...]

# Alternative approach using countByKey() action
word_count_dict = words_rdd.countByKey()  # This creates (word, 1) pairs
internally
print("Word count dictionary:", dict(word_count_dict))

# Sort by count (descending)
sorted_counts = word_counts.sortBy(lambda x: x[1], ascending=False)
print("Sorted by count:", sorted_counts.collect())
```

## Advanced RDD Operations

## Custom Partitioning

Understanding and controlling data partitioning is crucial for performance:

```python
# Check number of partitions
print("Default partitions:", numbers.getNumPartitions())

# Repartition data
repartitioned = numbers.repartition(4)
print("After repartitioning:", repartitioned.getNumPartitions())

# Coalesce - reduce partitions efficiently (only reduces, doesn't increase)
coalesced = numbers.coalesce(2)
print("After coalescing:", coalesced.getNumPartitions())
```

```python
# Custom partitioner for key-value RDDs
from pyspark import HashPartitioner

# Hash partitioning ensures keys with same hash go to same partition
partitioned_sales = sales_data.partitionBy(3, HashPartitioner(3))
print("Custom partitioned sales partitions:",
partitioned_sales.getNumPartitions())
```

## Caching and Persistence

Caching is crucial when you'll reuse an RDD multiple times:

```python
# Create an expensive computation
expensive_rdd = numbers.map(lambda x: x ** 2).filter(lambda x: x > 10)

# Cache the RDD in memory
expensive_rdd.cache()  # Same as expensive_rdd.persist(StorageLevel.MEMORY_ONLY)

# First action will compute and cache the RDD
result1 = expensive_rdd.collect()
print("First computation:", result1)

# Second action will use cached data (much faster)
result2 = expensive_rdd.count()
print("Count from cached data:", result2)

# Different storage levels
from pyspark import StorageLevel

# Persist with different storage levels
expensive_rdd.persist(StorageLevel.MEMORY_AND_DISK)  # Spill to disk if memory
full
expensive_rdd.persist(StorageLevel.MEMORY_ONLY_SER)  # Serialized in memory (space
efficient)
expensive_rdd.persist(StorageLevel.DISK_ONLY)       # Only on disk

# Unpersist when no longer needed
expensive_rdd.unpersist()
```

## Real-World Example: Log Analysis

Let's put everything together with a realistic log analysis scenario:

```python
# Simulate web server logs
log_data = [
    "192.168.1.1 - - [01/Jan/2023:10:00:00] GET /index.html 200 1024",
```

```python
        "192.168.1.2 - - [01/Jan/2023:10:01:00] GET /about.html 200 2048",
        "192.168.1.1 - - [01/Jan/2023:10:02:00] GET /contact.html 404 512",
        "192.168.1.3 - - [01/Jan/2023:10:03:00] POST /api/data 500 256",
        "192.168.1.2 - - [01/Jan/2023:10:04:00] GET /index.html 200 1024",
        "192.168.1.1 - - [01/Jan/2023:10:05:00] GET /api/users 200 4096"
]

logs_rdd = sc.parallelize(log_data)

# Parse log entries
def parse_log(log_line):
    parts = log_line.split()
    ip = parts[0]
    method = parts[5].strip('"')
    url = parts[6]
    status = int(parts[7])
    size = int(parts[8])
    return (ip, method, url, status, size)

parsed_logs = logs_rdd.map(parse_log)

# Analysis 1: Count requests per IP
ip_counts = parsed_logs.map(lambda x: (x[0], 1)).reduceByKey(lambda a, b: a + b)
print("Requests per IP:", ip_counts.collect())

# Analysis 2: Find all 404 errors
error_404 = parsed_logs.filter(lambda x: x[3] == 404)
print("404 errors:", error_404.collect())

# Analysis 3: Total bandwidth by status code
bandwidth_by_status = parsed_logs.map(lambda x: (x[3], x[4])).reduceByKey(lambda
a, b: a + b)
print("Bandwidth by status:", bandwidth_by_status.collect())

# Analysis 4: Most popular URLs
url_popularity = parsed_logs.map(lambda x: (x[2], 1)).reduceByKey(lambda a, b: a +
b)
top_urls = url_popularity.sortBy(lambda x: x[1], ascending=False).take(3)
print("Top URLs:", top_urls)
```

## Performance Considerations and Best Practices

Understanding how to write efficient RDD operations is crucial for production applications:

## Prefer reduceByKey over groupByKey

```python
# Inefficient approach
grouped = sales_data.groupByKey()
totals_inefficient = grouped.mapValues(sum)


# Efficient approach
totals_efficient = sales_data.reduceByKey(lambda x, y: x + y)
```

The reduceByKey approach is more efficient because it performs local aggregation before shuffling data across the network.

## Use appropriate actions

```python
# If you only need to check if data exists, don't use collect()
# Inefficient:
if len(numbers.collect()) > 0:
    print("Data exists")


# Efficient:
if not numbers.isEmpty():
    print("Data exists")
```

## Cache strategically

Only cache RDDs that you'll use multiple times, and remember to unpersist when done to free memory.

## Summary

RDD operations in PySpark provide a powerful abstraction for distributed data processing. The key concepts to remember are:

Transformations are lazy and create new RDDs without executing immediately. Actions trigger execution and return results to the driver program. Understanding the difference between transformations and actions helps you reason about when computations actually happen and how to optimize your Spark applications.

Key-value operations like reduceByKey, groupByKey, and join are fundamental for many data processing tasks. Proper partitioning and caching strategies can significantly improve performance for complex workflows.

The lazy evaluation model allows Spark to optimize the entire computation graph, but it also means you need to understand when your code actually executes. Actions like collect() can be expensive with large datasets, so use sampling and aggregation actions when possible.

Remember that RDDs are immutable and distributed, which enables fault tolerance but requires thinking differently about data processing compared to traditional imperative programming approaches.