

What are the different types of data structures in Spark

| Data Structure                        | Description  |
|---------------------------------------|--|
| Resilient Distributed Datasets (RDDs) | Immutable, distributed collections of data that can be operated on in parallel.                |
| DataFrames                            | Structured data that can be manipulated using SQL-like commands.                               |
| Datasets                              | Strongly typed data that can be operated on using either the functional or the imperative API. |
| Pair RDDs                             | RDDs that contain key-value pairs.   |
| Broadcast variables                   | Immutable variables that are shared across all workers in a Spark cluster.                     |
| Accumulators                          | Mutable variables that are used to track aggregate values across a Spark cluster.              |

Spark data structures can be operated on using a variety of different types of operations. **These operations can be classified into two main categories: transformations and actions.**

- **Transformations** are operations that create a new data structure from an existing one. They do not return any value, but they can be used to prepare the data for further processing. Some examples of transformations include:
  - map: Applies a function to each element in a dataset.
  - filter: Selects elements from a dataset that match a certain criteria.
  - reduce: Aggregates the elements in a dataset using a function.

- **Actions** are operations that return a value from a Spark job. They are typically used to store the results of a computation or to display the results to the user. Some examples of actions include:
  - `count`: Counts the number of elements in a dataset.
  - `sum`: Sums the values in a dataset.
  - `saveAsTextFile`: Saves the results of a computation to a file.

**Spark distributes data** across a cluster of machines using a process called partitioning. When a Spark job is submitted, the data is divided into a number of partitions, each of which is assigned to a different worker node in the cluster. The number of partitions is determined by the `spark.default.parallelism` property.

The partitioning of data in Spark is important for two reasons:

- It allows Spark to parallelize the processing of data.
- It allows Spark to keep the data in memory, which can improve performance.

There are two main types of partitioning in Spark: **hash partitioning** and **range partitioning**.

**Spark's job, stage, and task are three concepts that are used to describe the execution of a Spark application.**

- **Job** is the top-level unit of execution in Spark. A job is created when an action is invoked. Actions are operations that return a value from a Spark job. For example, the `count()` action returns the number of elements in a dataset.
- **Stage** is a sequence of tasks that can all be run together, in parallel, without a shuffle. Shuffles are operations that involve moving data between different partitions. For example, the `reduce()` operation requires a shuffle, because it aggregates the elements in a dataset using a function.
- **Task** is a single unit of work that is executed by a worker node in a Spark cluster. Tasks are typically responsible for performing a transformation or an action on a partition of data.