

How to Execute a Spark Program? 🚀

You can execute a Spark program using **multiple methods** depending on your setup and requirements.

1 Using spark-submit (Recommended for Production)

- `spark-submit` is the **official way** to run Spark applications in **cluster mode** or locally.
- Works for **both Python (PySpark) and Scala/Java**.
- Supports various **cluster managers** (Standalone, YARN, Mesos, Kubernetes).

💡 **Command:**

```
sh
CopyEdit
spark-submit --master local[*] my_spark_script.py
```

📌 **Options:**

- `--master local[*]` → Runs locally using all available CPU cores.
 - `--master yarn` → Runs on Hadoop YARN cluster.
 - `--master spark://<host>:<port>` → Runs on a Spark Standalone cluster.
-

2 Using PySpark Interactive Shell (For Testing & Debugging)

- Run Spark commands interactively in a Python shell.
- Best for **quick testing** and **data exploration**.

💡 **Command:**

```
sh
CopyEdit
pyspark
```

📌 **Then, run Spark commands:**

```
python
CopyEdit
```

```
df = spark.read.csv("data.csv", header=True, inferSchema=True)
df.show()
```

3 Running Spark in Jupyter Notebook (For Development & Testing)

- Use **Jupyter Notebook** for interactive data analysis with PySpark.
- Requires setting up `pyspark` and `findspark`.

💡 Setup Commands:

```
python
CopyEdit
import findspark
findspark.init()
from pyspark.sql import SparkSession
spark = SparkSession.builder.appName("Test").getOrCreate()
```

Then, you can run Spark queries in Jupyter Notebook.

4 Running Spark from an IDE (PyCharm, IntelliJ, VS Code)

- Useful for **development & debugging**.
- Requires **Spark configuration** inside the IDE.

💡 Example (Python - PyCharm/VS Code)

1. Set up **PySpark environment** (`pip install pyspark`).
2. Write a Python script (`my_spark_script.py`):

```
python
CopyEdit
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("Test").getOrCreate()
df = spark.read.csv("data.csv", header=True, inferSchema=True)
df.show()
```

3. Run the script inside the IDE.

For Java/Scala (IntelliJ)

- Set up **Spark dependencies** using **Maven** or **SBT**.
 - Write a Spark program and run it inside IntelliJ.
-

5 Running Spark in Databricks (For Cloud-based Execution)

- If using **Databricks**, upload the script as a **notebook**.
- Use `%python`, `%scala`, `%sql`, or `%r` magic commands.
- Run the notebook interactively in **Databricks clusters**.

Example:

```
python
CopyEdit
df = spark.read.csv("/mnt/data/sample.csv", header=True, inferSchema=True)
display(df)
```

6 Running Spark on a Hadoop YARN Cluster

- Requires a Hadoop cluster with **YARN Resource Manager**.
- Submit jobs using `spark-submit --master yarn`.

Command:

```
sh
CopyEdit
spark-submit --master yarn --deploy-mode cluster my_spark_script.py
```

Modes:

- `--deploy-mode client` → Runs driver on the local machine.
 - `--deploy-mode cluster` → Runs driver inside the cluster.
-

Summary of Different Execution Methods

Method	Use Case	Command/Setup
--------	----------	---------------

spark-submit	Production Deployment	<code>spark-submit my_script.py</code>
PySpark Shell	Quick Testing	<code>pyspark</code>
Jupyter Notebook	Interactive Development	<code>findspark + spark.read.csv()</code>
IDE (VS Code, PyCharm, IntelliJ)	Local Development	Write & Run Python/Scala scripts
Databricks	Cloud-Based Execution	Run Spark notebooks
YARN Cluster	Large-scale Distributed Processing	<code>spark-submit --master yarn</code>

How Spark Distributed Model Works?

Apache Spark follows a **distributed computing model** designed for **fault tolerance, parallel processing, and scalability**. It processes **huge datasets** efficiently across multiple nodes in a cluster.

1 Key Components of Spark’s Distributed Architecture

1. Driver Program

- The **entry point** of a Spark application.
- Runs on the **client machine** or **cluster** (depending on execution mode).
- Converts the user’s Spark code into a **Directed Acyclic Graph (DAG)**.
- Manages **task scheduling** and **result collection**.

 **Example:**

```
python
CopyEdit
spark = SparkSession.builder.appName("Test").getOrCreate()
```

This initializes the **Driver Program**.

2. Cluster Manager

- Allocates **resources (CPU, memory)** to Spark applications.
- Spark supports multiple **Cluster Managers**:
 - **Standalone Mode** → Spark's built-in cluster manager.
 - **YARN** → Runs Spark on a Hadoop cluster.
 - **Mesos** → General-purpose cluster manager.
 - **Kubernetes** → Runs Spark in containers.

Example:

```
sh
CopyEdit
spark-submit --master yarn my_script.py
```

This runs Spark on a YARN cluster.

3. Executors

- Executors are **worker processes** that **run tasks** and **store data**.
- Each Spark application has its own set of **dedicated Executors**.
- Executors **communicate with the Driver** and **execute transformations/actions**.
- Once an application completes, the Executors **shut down**.

Example of how Executors work:

1. Driver **assigns tasks** to Executors.
 2. Executors **process the data** in parallel.
 3. Results are **sent back** to the Driver.
-

4. Tasks & Stages Execution

- A **Task** is the smallest unit of execution (e.g., processing a partition).
- A **Stage** is a collection of tasks executed in parallel.
- Spark builds a **Directed Acyclic Graph (DAG)**:
 - **Narrow Transformations** (e.g., `map`, `filter`) → Do not cause data shuffling.
 - **Wide Transformations** (e.g., `groupBy`, `join`) → Require **shuffling**, meaning data is **moved between nodes**.

Example:

```
python
CopyEdit
rdd = spark.sparkContext.textFile("data.txt")
rdd2 = rdd.map(lambda x: x.upper()) # Narrow transformation
rdd3 = rdd2.groupBy(lambda x: x[0]) # Wide transformation (shuffle)
```

2 How Spark Distributes Data & Computation?

1 Dataset is divided into Partitions

- Spark **splits data into partitions** (logical chunks).
- Each partition is processed by **one Executor at a time**.
- Example: A **1TB dataset** can be **divided into 100 partitions**, each processed in parallel.

2 Driver creates RDDs (Resilient Distributed Datasets)

- RDDs are **immutable**, **fault-tolerant**, and **distributed** across nodes.
- Transformations create **new RDDs** without modifying existing ones.

3 Task Execution in Parallel

- Each partition is assigned a **task**, which runs on an **Executor**.
- Executors run tasks **in parallel** to improve performance.

4 Shuffling (if needed)

- If an operation **requires data movement** across nodes (e.g., `groupByKey()`), Spark performs **shuffling**.
- **Shuffling is expensive**, so Spark optimizes it using **broadcast joins, caching, and partitioning**.

5 Final Action Triggers Execution

- When an **action** like `.collect()`, `.count()`, or `.saveAsTextFile()` is called, Spark **executes all transformations** in the DAG.
-

3 Example: Spark Execution Flow

 Consider this PySpark code:

python
CopyEdit

```
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("Example").getOrCreate()
data = [("Alice", 25), ("Bob", 30), ("Charlie", 35)]
df = spark.createDataFrame(data, ["Name", "Age"])

df_filtered = df.filter(df.Age > 28)  # Lazy Transformation
df_filtered.show()  # Action (Execution starts here)
```

- ♦ **Step-by-step execution:**
- 1 **Driver Program** starts the application.
 - 2 **Cluster Manager** allocates resources.
 - 3 **RDD/DataFrame is partitioned** and distributed across Executors.
 - 4 **Transformations (filter)** are applied (but not executed yet).
 - 5 **Action (show())** triggers execution → DAG is executed.
 - 6 **Executors compute tasks in parallel** and return results.

4 Summary: Spark Distributed Model

Component	Function
Driver Program	Converts code into DAG and manages execution.
Cluster Manager	Allocates resources (Standalone, YARN, Mesos, Kubernetes).
Executors	Perform actual computations on worker nodes.
RDD/DataFrames	Distributed datasets stored across multiple nodes.
Tasks & Stages	Tasks execute transformations, stages optimize execution.

🔥 Spark enables high-speed distributed computing by breaking down tasks and running them in parallel! 🚀