



## **Chapter 2**

# **Downloading/Setting up Apache Spark**



# What is Homebrew?

- Homebrew is a popular package manager for macOS.
- It allows users to easily install, update, and manage software packages.

## Benefits from Homebrew:

- Streamlines software installation: No need to download and compile software manually.
- Offers a vast library of software packages for developers and users.
- Maintained by a community of contributors, ensuring up-to-date packages.

Basic commands: `brew install`, `brew update`, `brew upgrade`, and more.



# Step 1 - Downloading Apache Spark for Python

## 1. Install Homebrew

Simplify the installation and management of various software packages on your system.

*Run the following command in your terminal:*

```
/bin/bash -c "$(curl -fsSL https://raw.githubusercontent.com/Homebrew/install/HEAD/install.sh)"
```

## 2. Install Python with Homebrew

Ensure you have Python installed

*Run the following command:* **brew install python**

## 3. Install Apache Spark for Python

Enable integration of Apache Spark with Python

*Run the following command:* **pip install pyspark**

Update command: **pip install --upgrade pyspark**



## Step 2 - Test

**1. Run the following command on the terminal:**

```
1 pyspark
2 # Used to start an interactive PySpark shell, which is an interactive
3 # Python shell with PySpark functionality enabled.
```

**You should expect this:**

```

1 ~ pyspark
2 Python 3.10.7 (v3.10.7:6cc6b13308, Sep 5 2022, 14:02:52) [Clang 13.0.0 (clang-1300.0.29.30)] on darwin
3 Type "help", "copyright", "credits" or "license" for more information.
4 Setting default log level to "WARN".
5 To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
6 23/10/01 16:54:08 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
7 Welcome to
8
9      ____
10     /  _ \   _ __   ___
11    /  _ \  / __ \  / _ \
12   /  _ \ /  __/ /  __/
13  /  _ \/_  \_\  \_\_\
14
15  version 3.5.0
16
17 Using Python version 3.10.7 (v3.10.7:6cc6b13308, Sep 5 2022, 14:02:52)
18 Spark context Web UI available at http://geralds-mbp.fritz.box:4040
19 Spark context available as 'sc' (master = local[*], app id = local-1696172049014).
20 SparkSession available as 'spark'.

```

To exit the shell, type: **quit()**

## An example of using the Python Interpretive shell:

```

3 pyspark
4 Python 3.10.7 (v3.10.7:6cc6b13308, Sep 5 2022, 14:02:52) [Clang 13.0.0 (clang-1300.0.29.30)] on darwin
5 Type "help", "copyright", "credits" or "license" for more information.
6 Setting default log level to "WARN".
7 To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
8 23/10/01 16:54:08 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
9 Welcome to
10
11      /\_/\
12     /__\/
13
14     version 3.5.0
15
16 Using Python version 3.10.7 (v3.10.7:6cc6b13308, Sep 5 2022 14:02:52)
17 Spark context Web UI available at http://geralds-mbp.fritz.box:4040
18 Spark context available as 'sc' (master = local[*, app id = local-1696172049014]).
19 SparkSession available as 'spark'.
20 >>> strings = spark.read.text("README.md")
21 # The command reads the content of the "README.md" file into a PySpark DataFrame named "strings."
22
23 >>> strings.show(10, truncate=False)
24 # This command displays the first 10 rows of the "strings" DataFrame without truncating the text in each cell.
25
26 +-----+
27 |value|
28 +-----+
29 |# Apache Spark|
30 |
31 |Spark is a unified analytics engine for large-scale data processing. It provides
32 |high-level APIs in Scala, Java, Python, and R, and an optimized engine that
33 |supports general computation graphs for data analysis. It also supports a
34 |rich set of higher-level tools including Spark SQL for SQL and DataFrames,
35 |pandas API on Spark for pandas workloads, MLlib for machine learning, GraphX for graph processing,
36 |and Structured Streaming for stream processing.
37 |
38 |<https://spark.apache.org/>
39 +-----+
40
41 only showing top 10 rows
42
43 >>> strings.count()
44 125

```

# Important Concepts

## Application

It consists of a driver program and executors on the cluster.

## Job

Parallel computation triggered by Spark actions like `save()` and `collect()`.

## Task

A single unit of work or execution that will be sent to a Spark executor.

## SparkSession

Entry point for Spark API interaction in shell/applications.

## Stage

Jobs split into dependent stages for parallel execution.

# Transformations, Actions, and Lazy Evaluation

Spark operations on distributed data fall into two categories:

**Transformations** and **Actions**.

**Transformations** create a Spark DataFrame starting from another one, without modifying the original data.

**Example:** Operations like `select()` or `filter()` return a new DataFrame.

**\*\*Transformations are evaluated *lazily***

Meaning that results are not computed immediately, but they are recorded or remembered as a lineage.

**Actions** RDD actions are operations that return the raw values

**Example:** Spark actions include `count()` and `collect()`

Transformations	Actions
<code>orderBy()</code>	<code>show()</code>
<code>groupBy()</code>	<code>take()</code>
<code>filter()</code>	<code>count()</code>
<code>select()</code>	<code>collect()</code>
<code>join()</code>	<code>save()</code>

# Example of Transformation

```
1 >>> strings = spark.read.text("README.md")
2 >>> filtered = strings.filter(strings.value.contains("Spark"))
3 >>> filtered.count()
4 20
```

- It includes two transformations: `read()` (to read data) and `filter()` (to filter data).
- The key point is that nothing happens (in the terminal) when transformations are applied: execution is deferred until an action, such as `count()`, is called.
- In this example, the actual processing and filtering of data occur only when `filtered.count()` is executed in the Spark shell.



# Narrow and Wide Transformations

Transformations can be classified into two groups:

**Narrow** and **Wide** Transformations

**Narrow** transformations are those where each output can be computed from a single input.

They are *efficient* and *fast*.

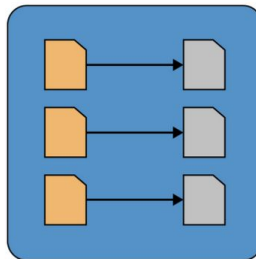
`filter()`, `map()`, `union()`, and `contains()`.

**Wide** transformations are those that require data from multiple input partitions to compute the output.

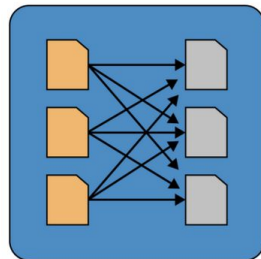
They are slower and less efficient as they require data shuffling to be performed.

`groupBy()`, `reduceByKey()`, `join()`, and `orderBy()`

Narrow Dependencies



Wide Dependencies



# The Spark UI

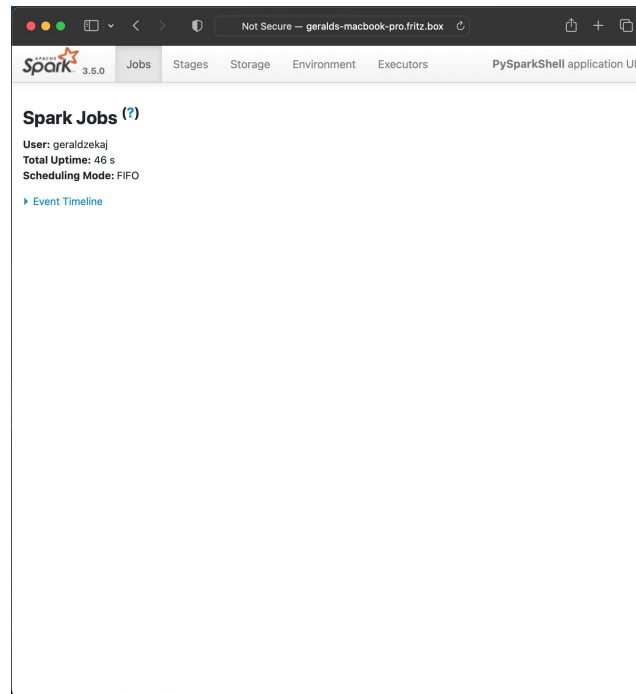
**Apache Spark** provides a graphical user interface for monitoring applications at various stages. It can be found on the first message when you launch SPARK

(<http://mf0023.homenet.telecomitalia.it:4040>).

The monitoring includes jobs, stages, and tasks.

## Key Information Available

- **Scheduler stages and tasks**
- **Memory Usage**
- **Environment Details**
- **Executor Information**
- **Spark SQL Queries**



# Code for Counting M&Ms

```
1 import sys
2 from pyspark.sql import SparkSession
3 from pyspark.sql.functions import count
4
5 # Check if this script is the main program
6 if __name__ == "__main__":
7     # Ensure that the correct number of command-line arguments are provided
8     if len(sys.argv) != 2:
9         print("Usage: python mnmcount.py <file>", file=sys.stderr)
10        sys.exit(-1)
11
12    # Create a SparkSession, which is the entry point to Spark functionality.
13    # If a SparkSession doesn't exist, it will create a new one.
14    spark = (SparkSession
15            .builder
16            .appName("PythonMnMCount") # Assign a name to the Spark application
17            .getOrCreate())
18
19    # Retrieve the M&M data set filename from the command-line arguments
20    mnm_file = sys.argv[1]
21
22    # Read the specified file into a Spark DataFrame. We use the CSV format
23    # and specify two options: "header" to indicate that the first row contains
24    # column names and "inferSchema" to automatically detect data types.
25    mnm_df = (spark.read.format("csv")
26            .option("header", "true")
27            .option("inferSchema", "true")
28            .load(mnm_file))
```

```
1 # We use the DataFrame high-level APIs to work with structured data. Note
2 # that we are exclusively using DataFrames and not RDDs (Resilient Distributed Datasets).
3 # In Spark, some functions return the same object, allowing us to chain function calls.
4
5 # Step 1: We select the fields "State," "Color," and "Count" from the DataFrame.
6 # Step 2: To group each state and its M&M color count, we use the groupBy() function.
7 # Step 3: We aggregate the counts of all colors and groupBy() "State" and "Color."
8 # Step 4: We order the results in descending order based on the "Total" count.
9 count_mnm_df = (mnm_df
10     .select("State", "Color", "Count")
11     .groupBy("State", "Color")
12     .agg(count("Count").alias("Total"))
13     .orderBy("Total", ascending=False))
14
15 # Show the resulting aggregations for all the states and colors;
16 # this displays the total count of each color per state.
17 # Note that "show()" is an action, which triggers the above query to be executed.
18 count_mnm_df.show(n=60, truncate=False)
19
20 # Print the total number of rows in the resulting DataFrame.
21 print("Total Rows = %d" % (count_mnm_df.count()))
```

```
1 # Here we focus on a single state.
2
3 # Step 1: Select all rows from the DataFrame.
4 # Step 2: Filter only those rows where the state is "CA."
5 # Step 3: Group the data by "State" and "Color" just like before.
6 # Step 4: Aggregate the counts for each color.
7 # Step 5: Order the results in descending order.
8
9 ca_count_mnm_df = (mnm_df
10     .select("State", "Color", "Count")
11     .where(mnm_df.State == "CA")
12     .groupBy("State", "Color")
13     .agg(count("Count").alias("Total"))
14     .orderBy("Total", ascending=False))
15
16 # Show the resulting aggregation for the state of California (CA).
17 # Similar to the previous example, the "show()" function is an action
18 # that will execute the computation and display the results.
19 ca_count_mnm_df.show(n=10, truncate=False)
20
21 # Stop the SparkSession.
22 spark.stop()
```

# Counting M&Ms for the Cookie Monster

1	+-----+-----+-----+
2	State Color  Count
3	+-----+-----+-----+
4	TX  Red  20
5	NV  Blue  66
6	CO  Blue  79
7	OR  Blue  71
8	WA  Yellow 93
9	+-----+-----+-----+
10	only showing top 5 rows

Fig. 1

1	+-----+-----+-----+
2	State Color  sum(Count)
3	+-----+-----+-----+
4	CA  Yellow 100956
5	CA  Brown  95762
6	CA  Green  93505
7	CA  Red  91527
8	CA  Orange 90311
9	CA  Blue  89123
10	+-----+-----+-----+

Fig. 3

1	+-----+-----+-----+
2	State Color  sum(Count)
3	+-----+-----+-----+
4	CA  Yellow 100956
5	WA  Green  96486
6	CA  Brown  95762
7	TX  Green  95753
8	TX  Red  95404
9	CO  Yellow 95038
10	NM  Red  94699
11	OR  Orange 94514
12	WY  Green  94339
13	NV  Orange 93929
14	TX  Yellow 93819
15	CO  Green  93724
16	CO  Brown  93692
17	CA  Green  93505
18	NM  Brown  93447
19	CO  Blue  93412
20	WA  Red  93332
21	WA  Brown  93082
22	WA  Yellow 92920
23	NM  Yellow 92747
24	NV  Brown  92478
25	TX  Orange 92315
26	AZ  Brown  92287
27	AZ  Green  91882
28	WY  Red  91768
29	AZ  Orange 91684
30	CA  Red  91527
31	WA  Orange 91521
32	NV  Yellow 91390
33	UT  Orange 91341
34	NV  Green  91331
35	NM  Orange 91251
36	NM  Green  91160
37	WY  Blue  91082
38	UT  Red  90995
39	CO  Orange 90971
40	AZ  Yellow 90946
41	TX  Brown  90736
42	OR  Blue  90526
43	CA  Orange 90311
44	OR  Red  90286
45	NM  Blue  90150
46	AZ  Red  90042
47	NV  Blue  90003
48	UT  Blue  89977
49	AZ  Blue  89971
50	WA  Blue  89886
51	OR  Green  89578
52	CO  Red  89465
53	NV  Red  89346
54	UT  Yellow 89264
55	OR  Brown  89136
56	CA  Blue  89123
57	UT  Brown  88973
58	TX  Blue  88466
59	UT  Green  88392
60	OR  Yellow 88129
61	WY  Orange 87956
62	WY  Yellow 87800
63	WY  Brown  86110
64	+-----+-----+-----+

Fig. 2