

Pyspark Introduction:

Overview of PySpark:

PySpark is the Python API for Apache Spark, a fast and general-purpose cluster computing system. Apache Spark provides in-memory data processing for big data workloads, making it efficient and suitable for a wide range of applications. PySpark allows developers to harness the power of Spark using Python, a language known for its simplicity and ease of use.

Code Example 1: Initializing PySpark and SparkContext

```
```python
from pyspark import SparkContext

Initialize a SparkContext
sc = SparkContext("local", "PySpark Example")

Perform a simple operation to check SparkContext
data = [18,25,14,7]
rdd = sc.parallelize(data)
result = rdd.map(lambda x: x *4).collect()

Print the result
print(result)

Stop the SparkContext when done
sc.stop()
```
```

In this example:

- We import the `SparkContext` class from the `pyspark` module.
- Initialize a `SparkContext` named "PySpark Example" with the master URL set to "local" (indicating that Spark will run in local mode on a single machine).
- Create a Resilient Distributed Dataset (RDD) from a list of numbers and perform a simple map operation to double each element.
- Collect the results and print them.
- Finally, stop the SparkContext to release resources.

Code Example 2: SparkSession for DataFrame API

```
```python
from pyspark.sql import SparkSession

Initialize a SparkSession
```

```

spark = SparkSession.builder.appName("PySpark DataFrame Example").getOrCreate()

Create a DataFrame from a list of tuples
data = [("John", 25), ("Jane", 30), ("Bob", 22)]
columns = ["Name", "Age"]
df = spark.createDataFrame(data, columns)

Show the DataFrame
df.show()

Perform a simple operation on the DataFrame
result_df = df.filter(df["Age"] > 25)
result_df.show()

Stop the SparkSession when done
spark.stop()
'''

```

In this example:

- We import the `SparkSession` class from the `pyspark.sql` module.
- Initialize a `SparkSession` with the application name "PySpark DataFrame Example."
- Create a DataFrame from a list of tuples and specify column names.
- Display the DataFrame using the `show()` method.
- Perform a simple filter operation on the DataFrame.
- Stop the SparkSession when the tasks are completed.

The `SparkSession` is the entry point to the DataFrame and SQL API in PySpark and provides a more user-friendly interface compared to the low-level RDD API.

Certainly! Let's delve deeper into `SparkContext` and `SparkSession` by exploring more functionality and use cases for each.

### Further Exploration of SparkContext:

**\*\*1. Using Broadcast Variables:\*\***

One powerful feature of Spark is the ability to use broadcast variables, which are read-only variables cached on each worker node rather than being sent over the network with tasks. This can greatly improve performance, especially when variables are large.

```

'''python
from pyspark import SparkContext

sc = SparkContext("local", "Broadcast Example")

```

```

Create a broadcast variable
broadcast_var = sc.broadcast([1, 2, 3, 4, 5])

Use the broadcast variable in a transformation
rdd = sc.parallelize([2, 3, 4])
result = rdd.map(lambda x: x * broadcast_var.value[0]).collect()

Print the result
print(result)

Stop the SparkContext
sc.stop()
```

```

In this example:

- We create a SparkContext as before.
- Create a broadcast variable named `broadcast_var` with a list of numbers.
- Use the broadcast variable in a map transformation on an RDD.
- The transformation multiplies each element of the RDD by the first element of the broadcast variable.
- Collect and print the results.

****2. Accumulators in SparkContext:****

Accumulators are variables that can be added to across multiple tasks in a parallel manner, making them useful for aggregating values across the workers.

```

```python
from pyspark import SparkContext

sc = SparkContext("local", "Accumulator Example")

Create an accumulator variable
accumulator_var = sc.accumulator(0)

Use the accumulator variable in a transformation
rdd = sc.parallelize([1, 2, 3, 4])
rdd.foreach(lambda x: accumulator_var.add(x))

Print the result
print("Accumulated Value:", accumulator_var.value)

Stop the SparkContext

```

```
sc.stop()
```
```

In this example:

- We create a SparkContext as before.
- Create an accumulator variable named `accumulator_var` initialized with 0.
- Use the accumulator variable in a `foreach` operation on an RDD to accumulate values.
- Print the final accumulated result.

Further Exploration of SparkSession:

****1. Reading and Writing Data using SparkSession:****

SparkSession provides a unified entry point for reading data from various sources and writing results to different sinks.

```
```python
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("Data IO Example").getOrCreate()

Reading data from a CSV file into a DataFrame
input_data = "path/to/input/data.csv"
df = spark.read.csv(input_data, header=True, inferSchema=True)

Show the DataFrame
df.show()

Writing DataFrame to Parquet format
output_data = "path/to/output/data.parquet"
df.write.parquet(output_data, mode="overwrite")

Stop the SparkSession
spark.stop()
```
```

In this example:

- We create a SparkSession as before.
- Read data from a CSV file into a DataFrame, specifying options such as headers and schema inference.
- Display the DataFrame using the `show()` method.
- Write the DataFrame to a Parquet file, specifying the output path and write mode.

****2. SQL Queries with SparkSession:****

SparkSession allows you to execute SQL queries on DataFrames, enabling the use of SQL for data manipulation.

```
```python
from pyspark.sql import SparkSession

spark = SparkSession.builder.appName("SQL Example").getOrCreate()

Creating a DataFrame
data = [("John", 25), ("Jane", 30), ("Bob", 22)]
columns = ["Name", "Age"]
df = spark.createDataFrame(data, columns)

Registering the DataFrame as a temporary SQL table
df.createOrReplaceTempView("people")

Executing a SQL query
result_df = spark.sql("SELECT * FROM people WHERE Age > 25")

Show the result DataFrame
result_df.show()

Stop the SparkSession
spark.stop()
```
```

In this example:

- We create a SparkSession as before.
- Create a DataFrame from a list of tuples.
- Register the DataFrame as a temporary SQL table named "people."
- Execute a SQL query on the registered table to filter records based on age.
- Display the result DataFrame.

These examples demonstrate more advanced features and use cases for both `SparkContext` and `SparkSession` in PySpark. They highlight the flexibility and power of PySpark in handling distributed data processing tasks.