Getting Started with PySpark DataFrames

Initializing a SparkSession:

Step 1: Import the required modules from pyspark.sql import SparkSession

Step 2: Create a SparkSession spark = SparkSession.builder \

.config("spark.some.config.option", "config-value") \ # Add any additional configurations if needed
.getOrCreate()

Step 3: Verify the SparkSession print(spark.version) # Print the Spark version to ensure the SparkSession is created

Loading data from CSV

Create a SparkSession
spark = SparkSession.builder \
.appName("Loading Data from CSV") \
.getOrCreate()

Load data from CSV file into a DataFrame csv_file_path = "path/to/your/csv/file.csv" df_csv = spark.read.csv(csv_file_path, header=True, inferSchema=True)

Show the first few rows of the DataFrame df csv.show()

Loading data from JSON

Load data from JSON file into a DataFrame json_file_path = "path/to/your/json/file.json"

df_json = spark.read.json(json_file_path) # Show the first few rows of the DataFrame df_json.show()

Loading data from Parquet

Load data from Parquet file into a DataFrame parquet_file_path = "path/to/your/parquet/file.parquet" df_parquet = spark.read.parquet(parquet_file_path)

Show the first few rows of the DataFrame df_parquet.show()

DataFrame Basics

In PySpark, a DataFrame is a distributed collection of data organized into named columns.For example, consider the following simple DataFrame representing sales data with distinct rows

+-----+

Creating DataFrames

PySpark's versatility extends beyond just working with structured data; it also offers the flexibility to create DataFrames from Resilient Distributed Datasets (RDDs) and various external data source

Creating DataFrames from RDDs

Let's explore how you can harness these capabilities.

Sample RDD rdd = spark.sparkContext.parllelize([(1, "Alice"), (2, "Bob"), (3, "Charlie")]) # Convert RDD to DataFrame

Show the DataFrame df_from_rdd.show ()

Working with Various Data Sources

Creating DataFrame from CSV csv_path = "path/to/data.csv" df_from_csv = spark.read.csv (csv_path, header=True, inferSchema-True)

json_path = "path/to/data.json" df_from_json = spark.read.json (json_path)

Creating DataFrame from JSON

Creating DataFrame from Parquet parquet path = "path/to/data.parquet" df_from_parquet = spark.read.parquet (parquet_path)

Basic DataFrame Operations

from pyspark.sql import SparkSession from pyspark.sql.functions import col # Create a SparkSession spark = SparkSession.builder \ .appName("DataFrame Operations Example") \ .getOrCreate()

Sample data for the DataFrame (1, "Alice", 28) (2, "Bob", 22),

(3, "Charlie", 35)

Define the DataFrame schema columns = ["ID", "Name", "Age"]

Create the DataFrame

Display DataFrame contents

Display the first 20 rows of the DataFrame df.show()

| Alice | 28

| Bob | 22

| Charlie | 35

+----+

df = spark.createDataFrame(data, columns)

| Name | Age

Check DataFrame schema

Check the DataFrame schema df.printSchema()

|-- Age: long (nullable = true)

|-- ID: long (nullable = true) |-- Name: string (nullable = true)

PySpark DataFrame Cheat Sheet

Selecting and Filtering Data

Original DataFrame:

| Alice | 28 | Bob | 22

Selecting specific columns from the DataFrame:

Select specific columns from the DataFrame selected columns = df.select("Name", "Age") # Show the DataFrame with selected columns print("DataFrame with Selected Columns:") selected columns.show()

Output:

DataFrame with Selected Columns:

+----+

Name | Age

Additionally, you can also use the col() function from pyspark.sql.functions to select columns. Here's how you can do it:

Using col() function to select columns selected_columns_v2 = df.select(col("Name"), col("Age"))

Show the DataFrame with selected columns using col() function print("DataFrame with Selected Columns (using col() function):") selected columns v2.show()

Output:

DataFrame with Selected Columns (using col() function): Name | Age

+-----

Filtering Data Elements

Creating a sample DataFrame to demonstrate each method:

Sample data for the DataFrame (1, "Alice", 28) (3, "Charlie", 35),

(5, "Eva", 25) # Define the DataFrame schema schema = StructType([

StructField("ID", IntegerType(), True). StructField("Name", StringType(), True), StructField("Age", IntegerType(), True)

Create the DataFrame df = spark.createDataFrame(data, schema)

Using filter() method

Method 1: Using filter() method filtered df = df.filter(df.Age > 25) # Show the filtered DataFrame

print("Filtered DataFrame (using filter()):")

filtered df.show()

Filtered DataFrame (using filter()):

+-----+ Alice Charlie 35 David 30

Using expr() Function

expr filtered df.show()

from pyspark.sql.functions import expr # Method 3: Using expr() function expr filtered df = df.filter(expr("Age > 25")) # Show the filtered DataFrame

print("Filtered DataFrame (using expr() function):")

Filtered DataFrame (using expr() function):

+----+ | Name | Age | Alice 28 | Charlie | 35 David 30

+----+

Using where() Method

Method 4: Using where() method where_filtered_df = df.where(df.Age > 25) # Show the filtered DataFrame print("Filtered DataFrame (using where() method):") where filtered df.show()

Filtering Data Elements

Filtered DataFrame (using where() method): +----+ | Name | Age | | Charlie | 3

Data Manipulation and Transformation

A sample DataFrame to demonstrate each operation:

| Name | Age Alice Bob | Charlie | 35

Adding new columns

Adding a new column "City" with a default value df_with_new_column = df.withColumn("City", lit("New York")) # Show the DataFrame with the new column print("DataFrame with New Column:") df with new column.show()

Output:

DataFrame with New Column:

Renaming Columns

df with renamed column.show()

Renaming the column "Name" to "Full Name" df with renamed column = df.withColumnRenamed("Name", "Full Name") # Show the DataFrame with the renamed column print("DataFrame with Renamed Column:")

DataFrame with Renamed Column: |Full Name | Age Bob

Dropping Columns # Dropping the column "Age" df dropped column = df.drop("Age") # Show the DataFrame with the dropped column

df dropped column.show()

Output: DataFrame with Dropped Column:

print("DataFrame with Dropped Column:")

Let's use a sample DataFrame to demonstrate some common data transformation examples:

from pyspark.sql import SparkSession from pyspark.sql.functions import col, lit, concat, expr # Create a SparkSession spark = SparkSession.builder \ .appName("Data Transformation") \ .getOrCreate() # Sample data for the DataFrame (3, "Charlie", 35)

Define the DataFrame schema columns = ["ID", "Name", "Age"]

Create the DataFrame

df = spark.createDataFrame(data, columns)

Using Functions to Transform Data

Using functions to transform data df with status = df.withColumn("Status", when(col("Age") > 25, "Adult").otherwise("Young")) # Show the DataFrame with the new "Status" column print("DataFrame with 'Status' Column:")

df with status.show()

DataFrame with Dropped Column:

Alice 28 Adult |Young |Adult | Bob | 22 | Charlie | 35

Data Manipulation and Transformation

Using Expressions to Transform Data

df_with_id_name = df.withColumn("ID_Name", concat(col("ID"), lit(" "), col("Name"))) # Show the DataFrame with the new "ID_Name" column print("DataFrame with 'ID_Name' Column:") df with id name.show()

Output:

DataFrame with 'ID Name' Column: Bob

Using SQL-like Expressions

Charlie 35

Using SQL-like expressions to transform data df_with_age_group = df.withColumn("Age_Group", expr("CASE WHEN Age <= 25 THEN 'Young' ELSE 'Adult' END")) # Show the DataFrame with the new "Age Group" column print("DataFrame with 'Age Group' Column:") df_with_age_group.show()

Output:

DataFrame with 'Age_Group' Column: | Age | Age_Groub |

|Adult |Young |Adult Bob

Aggregating and Grouping Data

from pyspark.sql import SparkSession from pyspark.sql.functions import col, sum, avg, count, min, max # Create a SparkSession spark = SparkSession.builder \ .appName("Aggregation Functions") \ # Sample data for the DataFrame

Define the DataFrame schema columns = ["Name", "Salary"] # Create the DataFrame

df = spark.createDataFrame(data, columns)

print("DataFrame:")

Show the DataFrame

DataFrame:

Aggregation using sum() Function # Using sum() function to calculate total salary total salary = df.select(sum("Salary")).collect()[0][0] print("Total Salary:", total_salary)

Total Salary: 750

Aggregation using avg() Function

Using avg() function to calculate average salary average_salary = df.select(avg("Salary")).collect()[0][0] print("Average Salary:", average_salary)

Average Salary: 150.0

Aggregation using count() Function

Using count() function to count the number of employees employee count = df.select(count("Name")).collect()[0][0] print("Number of Employees:", employee_count)

Number of Employees: 5

Aggregating and Grouping Data

Aggregation using min() and max() Functions

Using min() and max() functions to find minimum and maximum salary min_salary = df.select(min("Salary")).collect()[0][0] max_salary = df.select(max("Salary")).collect()[0][0] print("Minimum Salary:", min_salary) print("Maximum Salary:", max_salary)

Output:

Minimum Salary: 100

Maximum Salary: 200

Grouping Data Elements

from pyspark.sql import SparkSession from pyspark.sql.functions import col, avg # Create a SparkSession spark = SparkSession.builder .appName("Grouping Data") \ # Sample data for the DataFrame

("Alice", "Department A", 100) ("Bob", "Department B", 200), "Charlie", "Department A", 150) ("David", "Department C", 120) ("Eva", "Department B", 180)

Define the DataFrame schema columns = ["Name", "Department", "Salary"] # Create the DataFrame df = spark.createDataFrame(data, columns) # Show the DataFrame

Output:

print("DataFrame:")

df.show()

DataFrame: | Department A | 100 | Department B | 200 |Department B | 18

Grouping data based on "Department" and calculating average salary grouped_data = df.groupBy("Department").agg(avg("Salary").alias("AvgSalary")) # Show the grouped data print("Grouped Data:")

grouped_data.show() Output:

-----|------|Department B | 190.0 |Department C | 120.0 |Department A | 125.0

Joins and Combining DataFrames

A sample example to illustrate the different join types:

DataFrame 1: Name Alice DataFrame 1: l Role

left join = df1.join(df2, on="ID", how="left")

Manager

Now, let's demonstrate different join types:

1. Inner Join:

2. Outer Join:

3. Left Join:

inner_join = dfl.join(df2, on="ID", how="inner") inner_join.show() Bob 2021-12-20 2021-12-20 15:45:00 Charlie | 2022-02-28 | 2022-02-28 11:00:00 | outer join = dfl.join(df2, on="ID", how="outer") outer_join.show()

left_join.show() 4. Right Join: right_join = df1.join(df2, on="ID", how="right") right_join.show()

Let's now use a few examples to illustrate how to combine DataFrames using different join types:

Consider two sample DataFrames: a. DataFrame orders: +----+

b. DataFrame products: +-----+ |product| price | apple | 1.5 banana

banana

| grape | 2.0

Joins and Combining DataFrames

1. Inner Join: inner_join = orders.join(products, on="product", how="inner") +-----+----+ | banana | 102 | 2 | 0.75 | left_join = orders.join(products, on="product", how="left") left join.show() +-----+ +----- | banana | 102 | 2 | 0.75 | | orange | 103 | 5 | null | 3. Right Join: right_join = orders.join(products, on="product", how="right") right_join.show()

Handling Missing Data

outer_join = orders.join(products, on="product", how="outer")

Identifying and handling missing data is a crucial step in data preprocessing and analysis. PySpark provides several techniques to identify and handle missing data in DataFrames. Let's explore these techniques:

Identifying Missing Data

grape | null | null | 2.0

+-----+----+

4. Outer Join:

PySpark represents missing data as null values. You can use various methods to identify missing data in DataFrames: isNull() and isNotNull()

Handling Missing Data Dropping Rows with Null Values

df _without _null = df.dropna()

 Filling Null Values df filled mean = df.fillna({ 'age': df.select(avg('age')).first() [0] })

from pyspark.ml.feature import Imputer

Imputation

imputer = Imputer(inputCols=['age'], outputCols=['imputed age']) imputed df = imputer.fit(df).transform(df) Adding an Indicator Column

df with indicator = df.withColumn('age missing', df['age'].isNull())

 Handling Categorical Missing Data df_filled_category = df.fillna({'gender': 'unknown'})

Working with Dates and Timestamps

Sample DataFrame with date and timestamp data as mentioned below DataFrame:

Current Date and Timestamp

print("DataFrame with Current Timestamp:")

from pyspark.sql.functions import current_date, current_timestamp # Adding columns for current date and timestamp df with current date = df.withColumn("CurrentDate", current date()) df_with_current_timestamp = df.withColumn("CurrentTimestamp", current_timestamp()) # Show the DataFrames print("DataFrame with Current Date:")

Date Difference

date diff df.show()

df with current date.show()

df_with_current_timestamp.show()

from pyspark.sql.functions import datediff date_diff_df = df.withColumn("DaysSince", datediff(current_date(), col("Date"))) # Show the DataFrame with date difference print("DataFrame with Date Difference:")

from pyspark.sql.functions import date_add, date_sub date add df = df.withColumn("DatePlus10Days", date add(col("Date"), 10) date sub df = df.withColumn("DateMinus5Days", date sub(col("Date"), 5)) # Show the DataFrames with date addition and subtraction print("DataFrame with Date Addition:") date_add_df.show() print("DataFrame with Date Subtraction:") **Advanced DataFrame Operations Window Functions**

Months Between

from pyspark.sql.functions import months_between

Show the DataFrame with months between

print("DataFrame with Months Between:")

Date Addition and Subtraction

months between df.show()

Suppose you have a DataFrame with sales data and you want to calculate the rolling average of sales for each product over a specific window size. from pyspark.sql import SparkSession from pyspark.sql.window import Window from pyspark.sql.functions import col, avg # Create a SparkSession spark = SparkSession.builder \ .appName("Window Functions") \ .getOrCreate() # Sample data for the DataFrame ("ProductA", "2022-01-01", 100), ("ProductA", "2022-01-02", 150), ("ProductA", "2022-01-03", 200), ("ProductA", "2022-01-04", 120) ("ProductA", "2022-01-05", 180) ("ProductB", "2022-01-01", 50),

window_spec = Window.partitionBy("Product").orderBy("Date").rowsBetween(-1, 1)

Working with Dates and Timestamps

Define the DataFrame schema columns = ["Product", "Date", "Sales"] # Create the DataFrame df = spark.createDataFrame(data, columns)

Calculate rolling average using window function df_with_rolling_avg = df.withColumn("RollingAvg", avg(col("Sales")).over(window_spec)) # Show the DataFrame with rolling average

Pivot Tables Suppose you have a DataFrame with sales data and you want to create a pivot table showing total sales for each product in different months.

from pyspark.sql import SparkSession from pyspark.sql.functions import col from pyspark.sql.pivot import PivotTable # Create a SparkSession spark = SparkSession.builder \
.appName("Pivot Tables") \ .getOrCreate() # Sample data for the DataFrame

("ProductB", "2022-01-02", 70

Define the window specification

df_with_rolling_avg.show()

("ProductA", "2022-01-01", 100). ("ProductA", "2022-02-01", 150) ("ProductA", "2022-01-01", 200) ("ProductB", "2022-02-01", 120) ("ProductB", "2022-02-01", 180)

Show the pivot table

pivot table.show()

Define the DataFrame schema

columns = ["Product", "Date", "Sales"] # Create the DataFrame df = spark.createDataFrame(data, columns) # Create a pivot table

pivot_table = df.groupBy("Product").pivot("Date").sum("Sales")

PySpark SQL Cheat Sheet: SQL Functions for

DataFrames

Filtering Data

from pyspark.sql import SparkSession from pyspark.sql.functions import col # Create a SparkSession spark = SparkSession.builder \ .appName("PySpark SQL Functions") \ .getOrCreate() # Sample data for the DataFrame

Define the DataFrame schema columns = ["Name", "Age"] # Create the DataFrame

df = spark.createDataFrame(data, columns)

("Charlie", 35)

Using SQL-like functions to filter data filtered data = df.filter(col("Age") > 25) selected columns = df.select("Name", "Age") # Show the filtered and selected data print("Filtered Data:") filtered data.show()

print("Selected Columns:")

selected columns.show()

Aggregation

from pyspark.sql.functions import avg, max # Using SQL-like functions for aggregation grouped_data = df.groupBy("Age").agg(avg("Age"), max("Age")) # Show the aggregated data print("Aggregated Data:") grouped_data.show()

PySpark SQL Cheat Sheet: SQL Functions for

DataFrames

Sorting Data

from pyspark.sql.functions import desc # Using SQL-like functions for sorting

desc_sorted_data = df.sort(col("Age").desc()) # Show the sorted data

sorted data = df.orderBy("Age")

print("Descending Sorted Data:") desc sorted data.show()

("Bob", "San Francisco"),

print("Sorted Data:")

sorted data.show()

Joining DataFrames # Sample data for another DataFrame ("Alice", "New York"),

("Eva", "Los Angeles") columns2 = ["Name", "City"] df2 = spark.createDataFrame(data2, columns2)

Using SQL-like functions to join DataFrames

ioined data = df.join(df2, on="Name", how="inner")

Show the joined data print("Joined Data:") joined data.show()

Performance Optimization Tips

efficient data processing. Here are some tips and best practices to consider:

DataFrames that are used in multiple transformations or actions.

Optimizing PySpark DataFrame operations is essential to achieve better performance and

1. Use Lazy Evaluation: PySpark uses lazy evaluation, which means that transformations on DataFrames are not executed immediately but are queued up. This allows PySpark to optimize and optimize the execution plan before actually performing computations.

2. Caching: Caching involves storing a DataFrame or RDD in memory so that it can be

3. Partitioning: Partitioning involves dividing your data into smaller subsets (partitions) based on certain criteria. This can significantly improve query performance as it reduces the amount of data that needs to be processed. Use .repartition() or .coalesce() to manage

4. Broadcasting: Broadcasting is a technique where smaller DataFrames are distributed to

lead to memory issues. Instead, try to perform operations using distributed computations.

6. Use Built-in Functions: Whenever possible, use built-in functions from the pyspark.sql.

reused efficiently across multiple operations. Use .cache() or .persist() to cache intermediate

worker nodes and cached in memory for join operations. This is particularly useful when you have a small DataFrame that needs to be joined with a larger one. 5. Avoid Using .collect(): Using .collect() brings all the data to the driver node, which can

functions module. These functions are optimized for distributed processing and provide better performance than custom Python functions. 7. Avoid Shuffling: Shuffling is an expensive operation that involves data movement

operations that require data to be rearranged.

and identify potential optimization opportunities.

cluster can significantly impact performance..

8. Optimize Joins: Joins can be performance-intensive. Try to avoid shuffling during joins by ensuring both DataFrames are appropriately partitioned and using the appropriate join strategy (broadcast, sortMerge, etc.).

between partitions. Minimize shuffling by using appropriate partitioning and avoiding

10. Hardware Considerations: Consider the cluster configuration, hardware resources, and the amount of data being processed. Properly allocating resources and scaling up the

9. Use explain(): Use the explain() method on DataFrames to understand the execution plan

11. Monitor Resource Usage: Keep an eye on resource usage, including CPU, memory, and

disk I/O. Monitoring can help identify performance bottlenecks and resource constraints... 12. Use Parquet Format: Parquet is a columnar storage format that is highly efficient for both reading and writing. Consider using Parquet for storage as it can improve read and write performance.

