

### ### 27. Integration with Pandas:

PySpark provides a convenient way to convert between PySpark DataFrames and Pandas DataFrames. This is useful when you want to leverage the capabilities of both PySpark and Pandas in your data processing tasks.

**\*\*Sample PySpark Code:\*\***

```
```python
from pyspark.sql import SparkSession
import pandas as pd

# Create a Spark session
spark = SparkSession.builder.appName("pyspark_pandas_integration").getOrCreate()

# Create a PySpark DataFrame
data = [("John", 28), ("Alice", 35), ("Bob", 40)]
columns = ["Name", "Age"]
pyspark_df = spark.createDataFrame(data, columns)

# Convert PySpark DataFrame to Pandas DataFrame
pandas_df = pyspark_df.toPandas()

# Display PySpark DataFrame
print("PySpark DataFrame:")
pyspark_df.show()

# Display Pandas DataFrame
print("\nPandas DataFrame:")
print(pandas_df)

# Convert Pandas DataFrame to PySpark DataFrame
new_pyspark_df = spark.createDataFrame(pandas_df)

# Display the new PySpark DataFrame
print("\nNew PySpark DataFrame:")
new_pyspark_df.show()

# Stop the Spark session
spark.stop()
```
```

In this example, we create a PySpark DataFrame, convert it to a Pandas DataFrame, and then convert the Pandas DataFrame back to a PySpark DataFrame.

### ### 28. Data Skewing:

Data skewness in Spark jobs can occur when certain keys have significantly more data than others, leading to performance issues. Identifying and addressing skewed data is crucial for optimizing Spark jobs.

**\*\*Sample PySpark Code:\*\***

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

# Create a Spark session
spark = SparkSession.builder.appName("data_skewing_example").getOrCreate()

# Create a sample DataFrame with skewed data
data = [("A", 100), ("B", 150), ("C", 50), ("D", 200)]
columns = ["Key", "Value"]
skewed_df = spark.createDataFrame(data, columns)

# Introduce skewness by duplicating data for key "D"
skewed_df = skewed_df.union(spark.createDataFrame([("D", 200)] * 800, columns))

# Perform an operation that can be affected by skewness
result_df = skewed_df.groupBy("Key").agg({"Value": "sum"})

# Display the result DataFrame
print("Skewed Result DataFrame:")
result_df.show()

# Addressing skewness by repartitioning the DataFrame
repartitioned_df = skewed_df.repartition("Key")

# Perform the operation on the repartitioned DataFrame
result_repartitioned_df = repartitioned_df.groupBy("Key").agg({"Value": "sum"})

# Display the result of the repartitioned DataFrame
print("\nResult after Repartitioning:")
result_repartitioned_df.show()

# Stop the Spark session
```

```
spark.stop()  
``
```

In this example, we create a `DataFrame` with skewed data and then address the skewness issue by repartitioning the `DataFrame` before performing the aggregation operation. Repartitioning helps distribute the data more evenly across partitions, reducing the impact of skewness on performance.