## Challenge 1: Write a PySpark function to remove

 duplicate rows from a DataFrame based on specific columns.

▼ Function: remove\_duplicates\_by\_columns

```
1 from pyspark.sql import DataFrame
 3 def remove_duplicates_by_columns(df: DataFrame, columns: list) -> DataFrame:
 4
 5
       Removes duplicate rows based on specified columns.
 6
 7
       Parameters:
8
       df (DataFrame): The input PySpark DataFrame.
       columns (list): A list of column names to check for duplicates.
9
10
11
       Returns:
12
      DataFrame: A new DataFrame with duplicates removed based on specified cc
13
       return df.dropDuplicates(columns)
14
15
```

Example Usage

```
1 from pyspark.sql import SparkSession
 3 spark = SparkSession.builder.appName("RemoveDuplicatesExample").getOrCreate(
 5 # Sample data
 6 data = [
      ("Alice", "HR", 1000),
      ("Bob", "IT", 1200),
8
      ("Alice", "HR", 1000), # duplicate
      ("Alice", "IT", 1000), # not a duplicate
10
11 ]
12
13 columns = ["name", "department", "salary"]
14 df = spark.createDataFrame(data, columns)
15
16 # Remove duplicates based on 'name' and 'department'
17 cleaned_df = remove_duplicates_by_columns(df, ["name", "department"])
18
19 cleaned_df.show()
20
```

## Output

```
1 +----+
2 | name|department|salary|
3 +----+
4 |Alice| HR| 1000|
5 | Bob| IT| 1200|
6 |Alice| IT| 1000|
7 +----+
```

The row ("Alice", "HR", 1000) appeared twice, but only one instance is kept based on the ["name", "department"] columns.

Challenge 2: Create a PySpark pipeline to read a CSV
file, filter out rows with null values, and write the result to a Parquet file.

- Objective:
- Read a CSV file.
- Filter out rows with any null values.
- Write the cleaned data to a Parquet file.

```
1 from pyspark.sql import SparkSession
3 def csv_to_parquet_pipeline(input_csv_path: str, output_parquet_path: str):
4
      # Step 1: Create Spark session
5
      spark = SparkSession.builder \
           appName("CSVtoParquetPipeline") \
 6
 7
           .get0rCreate()
 8
      # Step 2: Read CSV file
9
      df = spark.read.option("header", "true").csv(input_csv_path)
10
11
12
      # Step 3: Filter out rows with any null values
      cleaned_df = df.dropna()
13
14
15
      # Step 4: Write to Parquet
      cleaned_df.write.mode("overwrite").parquet(output_parquet_path)
16
17
      print(f"√ Pipeline completed. Parquet written to: {output_parquet_patl
18
19
20
      # Optional: Stop Spark session
21
      spark.stop()
22
```

## Example Usage

1 csv\_to\_parquet\_pipeline("/content/drive/MyDrive/pyspark csv data for practi

## Challenge 3: Implement a window function to rank salespeople based on total sales by region.

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.types import StringType, StructType, IntegerType, StructFi
3 from pyspark.sql.functions import col, sum, desc, rank
4 from pyspark.sql.window import Window
```

```
5
6 #Step:1 - Create spark session
7 spark = SparkSession.builder.appName("rankSalesPeopleByRegion").getOrCreate
9 schema = StructType([
                         StructField('salesperson', StringType(), True),
                         StructField('region', StringType(), True),
11
12
                         StructField('sale_amount', IntegerType(), True)
13
                       1)
14
15 #Step:2 - Read CSV file & create a df
16 df = spark.read.option('header', 'true').csv("/content/drive/MyDrive/pyspar
17
                                                       schema=schema)
18 # sales_df = df.groupBy('salesperson', 'region').agg(sum('sale_amount')).al
19 # df.show()
20 # df.printSchema()
21
22 df = df.withColumn('sale_amount', col('sale_amount').cast('int'))
23
24 # Step:3 - create a new df
25 sales_df = df.groupBy('salesperson', 'region').agg(sum('sale_amount').alias
26 # sales df.show()
27
28 #Step:4 - Apply a window function
29 window_spec = Window.partitionBy('region').orderBy(desc('total_sales'))
30
31 ranked_df = sales_df.withColumn('Rank', rank().over(window_spec))
33 ranked df.show()
```



+		
salesperson	region	  total_sales
Ivan	South	1350
David	South	1100
Judy	East	1000
Alice	East	1200
Bob	West	800
Eve	West	1500
Charlie	East	950
Frank	East	700
Heidi	West	500
Grace	South	1250
		L

<b>L</b>				
salesperson	region	total_sale	s	Rank
Alice	East	120	00	1
Judy	East	100	)0 j	2 į
Charlie	East	95	50 İ	3
Frank	East	70	)0 j	4
Ivan	South	135	i0	1
Grace	South	125	0	2
David	South	110	00	3
Eve	West	150	00	1
Bob	West	80	00	2
Heidi	West	50	00	3
+		<b></b>	+	+

Challenge 4: Write a PySpark SQL query to calculate

 the average salary by department, including only employees with more than 3 years of experience.



Write a PySpark SQL query to:

Filter employees with more than 3 years of experience

Group by department

Calculate the average salary

```
1 from pyspark.sql import SparkSession
 2 from pyspark.sql.types import StructField, StructType, StringType, Integerl
 3 from pyspark.sql.functions import avg
 5 # Step-1: create a spark session
 6 spark = SparkSession.builder.appName("AvgSal_Deptwise").getOrCreate()
 8 schema = StructType([
       StructField('employee_id', IntegerType(), True),
       StructField('name', StringType(), True),
10
       StructField('department', StringType(), True),
11
       StructField('salary', IntegerType(), True),
12
       StructField('experience', IntegerType(), True)
13
14 1)
15
16 # Step-2 : Read csv
17 df = spark.read.option('header', 'true').csv('/content/drive/MyDrive/pyspar
                                                        schema=schema)
18
19 # df.show()
20 # df.printSchema()
21
22 # Step-3 : Filter employees with more than 3 years of experience
23 df = df.filter(df['experience'] > 3)
24 # df.show()
25
26 # Step-4 : Calculate the average salary department-wise
27 df = df.groupBy('department').agg(avg('salary').alias('avg_sal_deptwise'))
28 df.show()
29
\rightarrow
    |department|avg_sal_deptwise|
             HR I
                          61000.0|
        Finance
                         90000.01
                         80500.0
```

Challenge 5: Implement a PySpark function to split a
 large DataFrame into smaller DataFrames based on a specific column value.

<sup>1</sup> from nychark cal impart CharkCoccion

```
T LLOW hAzhark*zdr TWholr Sharksezzton
 2
 3 spark = SparkSession.builder.appName("Split_LargeDF_to_SmallerDFs").getOrCre
 5 schema = StructType([
      StructField('employee_id', IntegerType(), True),
 6
 7
      StructField('name', StringType(), True),
      StructField('region', StringType(), True),
8
9
      StructField('department', StringType(), True),
10
      StructField('salary', IntegerType(), True)
11 1)
12
13 # Step-1: Read CSV
14 df = spark.read.option('header', 'true').csv('/content/drive/MyDrive/pyspark
15
                                                 , schema = schema)
16 df.show()
17
18 # Step-2: Let's split the data on Region, let's store unique regions in regi
19 regions = [row['region'] for row in df.select('region').distinct().collect()
20
21 # Step-3: Create a dictionary of data frames for each region
22 dfs_by_region = {region:df.filter(df.region == region) for region in regions
23
24 dfs_by_region
25 print(dfs_by_region)
26
27 dfs_by_region['East'].show()
28 dfs_by_region['West'].show()
29 dfs_by_region['North'].show()
30 dfs_by_region['South'].show()
```



			L	L
employee_id	name	region	department	salary
1	Bob  Charlie   David   Eve   Frank   Grace   Heidi	East West North East West South East North	IT IT   Finance   HR IT   Finance   HR	58000    79000    91000    62000
10		West		
+				<del>-</del>

{'South': DataFrame[employee\_id: int, name: string, region: string, departm
+-----+
|employee\_id| name|region|department|salary|
+-----+
1	Alice	East	HR	60000
4	David	East	Finance	90000
7	Grace	East	Finance	91000