

## Assignment 1: Working with CSV Data (employee\_data.csv)

### Tasks:

#### 1. Load the CSV data:

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import col, to_date, avg, count
3
4 # Create a SparkSession
5 spark = SparkSession.builder.appName("EmployeeDataAnalysis").getOrCreate()
6 # Convert into dbfs
7 dbutils.fs.cp("file:/Workspace/Users/azuser2130_mml.local@techademy.com/Employees_data.csv", "dbfs:/FileStore/Employees_data.csv")
8 # Load CSV data
9 df = spark.read.format("csv").option("header", "true").load(["dbfs:/FileStore/Employees_data.csv"])
10 # Display the first 10 rows
11 df.show(10)
```

Output Terminal Debug Console

df: pyspark.sql.dataframe.DataFrame = [EmployeeID: string, Name: string ... 3 more fields]

EmployeeID	Name	Department	JoiningDate	Salary
1001	John Doe	HR	2021-01-15	55000
1002	Jane Smith	IT	2020-03-10	62000
1003	Emily Johnson	Finance	2019-07-01	70000
1004	Michael Brown	HR	2018-12-22	54000
1005	David Wilson	IT	2021-06-25	58000
1006	Linda Davis	Finance	2020-11-15	67000
1007	James Miller	IT	2019-08-14	65000
1008	Barbara Moore	HR	2021-03-29	53000

#### 2. Data Cleaning:

```
# Remove rows where Salary is less than 55,000
df_cleaned = df.filter(col('Salary') >= 55000)
df_cleaned.show()

# Filter employees who joined after 2020
df_cleaned = df_cleaned.filter(to_date(df_cleaned['JoiningDate']) >= '2020-01-01')
df_cleaned.show()
```

(2) Spark Jobs

df\_cleaned: pyspark.sql.dataframe.DataFrame = [EmployeeID: string, Name: string ... 3 more fields]

EmployeeID	Name	Department	JoiningDate	Salary
1001	John Doe	HR	2021-01-15	55000
1002	Jane Smith	IT	2020-03-10	62000
1003	Emily Johnson	Finance	2019-07-01	70000
1005	David Wilson	IT	2021-06-25	58000
1006	Linda Davis	Finance	2020-11-15	67000
1007	James Miller	IT	2019-08-14	65000

EmployeeID	Name	Department	JoiningDate	Salary
1001	John Doe	HR	2021-01-15	55000
1002	Jane Smith	IT	2020-03-10	62000
1005	David Wilson	IT	2021-06-25	58000
1006	Linda Davis	Finance	2020-11-15	67000

### 3. Data Aggregation:

▶

✓ Just now (1s)

6

```
1 # Find the average salary by Department
2 average_salary_by_department = df_cleaned.groupby('Department').agg(avg('Salary').alias('AverageSalary'))
3 print("Average Salary by Department:\n")
4 average_salary_by_department.show()
5
6 # Count the number of employees in each Department
7 employee_count_by_department = df_cleaned.groupby('Department').count()
8 print("Employee Count by Department:\n")
9 employee_count_by_department.show()
10
```

Output

Terminal

Debug Console

Average Salary by Department:

Department	AverageSalary
HR	55000.0
Finance	67000.0
IT	60000.0

Employee Count by Department:

Department	count
HR	1
Finance	1
IT	2

### 4. Write the Data to CSV:

Save the cleaned data (from the previous steps) to a new CSV file.

▶

✓ Just now (1s)

6

```
# Save the cleaned data to a new CSV file
df_cleaned.write.mode("overwrite").csv("dbfs:/path/to/cleaned_employee_data.csv", header=True)
```

▶ (1) Spark Jobs