Coding Challenge PySpark

Name: Divya Sree Murali Date: 26/11/2024

1.Explain ETL (Extract, Transform, Load) with PySpark(in your own words)

ETL (Extract, Transfer, Load) is a process used to manage data workflows. It is used to Combine data from multiple sources into a unified view and format data for better analysis. PySpark is the Python API for Apache Spark that provides a scalable and efficient way to implement ETL. By leveraging PySpark, the ETL process becomes more efficient and scalable.

1. Extract:

This phase involves fetching of raw data from various sources like relational databases, files, streaming data etc. PySpark supports multiple formats like CSV, JSON files. To read the data, PySpark uses SparkSession.read keyword

2. Transform:

The raw data extracted often requires cleaning, enrichment, or reshaping to meet business requirements. PySpark provides tools to perform these transformations at scale:

- Use Data Frame API for operations like filtering, joining, grouping, or aggregating.
- Use PySpark SQL for querying data with SQL-like syntax.
- Apply functions using with Column or select to modify or add new columns.

3. Load

The processed data is then written to a target location, such as a data warehouse, database, or storage system. PySpark supports writing data in various formats and to different storage backends and Formats supported include CSV, Parquet, ORC, JSON, etc.

The Advantages of using PySpark for ETL is that:

PySpark is highly scalable such that it can handle large datasets distributed across clusters. It is flexible and hence supports multiple data sources and formats. It is optimized speed for parallel processing. It is versatile and can easily integrate with numerous data sources and works well with other Big Data Tools

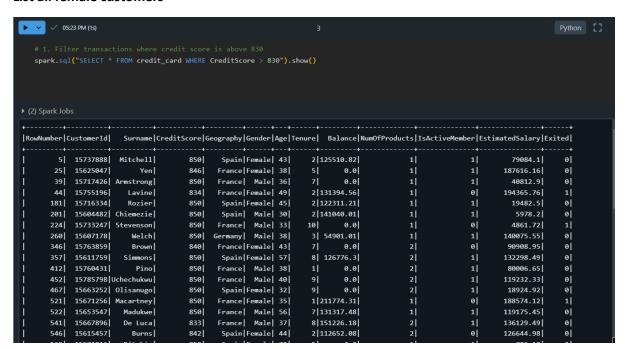
2.. Using spark sql Transformations such as Filter, Join, Simple Aggregations, GroupBy on the case study dataset

Using spark sql

1.Filter

List all customers having credit score greater tan 830

List all female customers



```
▶ ∨ √ 05:24 PM (<1s)
                                                                                                                               Python []
  spark.sql("SELECT * FROM credit card WHERE Gender = 'Female'").show()
|RowNumber|CustomerId| Surname|CreditScore|Geography|Gender|Age|Tenure| Balance|NumOfProducts|IsActiveMember|EstimatedSalary|Exited|
        1 15634602
                                       619
                                             France|Female| 42|
                                                                                                                  101348.88
                      Hargrave
           15647311
                          Hill|
                                       608
                                              Spain|Female| 41|
                                                                    1| 83807.86|
                                                                                                                  112542.58
           15619304
                          Onio
                                       502
                                              France|Female| 42|
                                                                    8 | 159660.8
                                                                                                          0|
                                                                                                                  113931.57
                                             France|Female| 39|
Spain|Female| 43|
           15701354
                          Bonil
                                       699
                                                                           0.01
                                                                                                          0|
                                                                                                                  93826.63
                                                                                                                                 0|
                                                                    2|125510.82|
           15737888
                      Mitchell|
                                                                                                          11
                                                                                                                                 0|
                                       850 l
                                                                                                                   79084.1
                                            Germany|Female| 29|
                                                                                                                  119346.88
        8 15656148
                                                                    4|115046.74|
                                                                                                                                 1|
                        Obinna|
                                       376
                                                                                                          0|
                                             France|Female| 34|
                                                                                                                  26260.98
       13 15632264
                          Kay
                                       476
                                                                   10
                                                                           0.0
                                                                                                          0|
                                                                                                                                 0|
       14 | 15691483 |
                          Chin
                                       549
                                              France|Female| 25|
                                                                            0.0
                                                                                                                  190857.79
                                                                                                          0|
                                                                                                                                 0|
       15 | 15600882
                                              Spain|Female| 35|
                                                                                                                  65951.65
                          Scott|
                                       635
                                                                            0.0
                                                                                                                                 0|
            15788218| Henderson|
                                       549
                                               Spain|Female| 24|
                                                                            0.0
                                                                                                                   14406.41
                                                                                                                                 0|
       20|
            15568982
                          Hao
                                       726
                                              France|Female| 24|
                                                                            0.0
                                                                                                                   54724.03
                                                                                                                                 0|
       22
            15597945| Dellucci|
                                       636
                                               Spain|Female| 32|
                                                                            0.0
                                                                                                                  138555.46
       23
           15699309 Gerasimov
                                       510
                                              Spain|Female| 38|
                                                                            0.0
                                                                                                          0|
                                                                                                                  118913.53
                                             France|Female| 38|
       25
           15625047
                          Yen
                                       8461
                                                                            0.01
                                                                                                          1|
                                                                                                                  187616.16
                                                                                                                                 0|
       29 | 15728693 | McWilliams |
                                       574 Germany Female 43
                                                                    3 | 141349.43 |
                                                                                            1
                                                                                                                  100187.43
                                                                                                                                 øl
            15589475
```

2.Joins

Inner Join

```
""").show()
|RowNumber|CustomerId| Surname|CreditScore|Geography|Gender|Age|Tenure| Balance|NumOfProducts|IsActiveMember|EstimatedSalary|Exited|
        1 15634602 Hargrave
                                      619 l
                                             France|Female| 42|
                                                                           0.01
                                                                                                           11
                                                                                                                    101348.88
                                                                                                                                  11
           15647311
                         Hi111
                                      608 l
                                              Spain|Female| 41|
                                                                    1 83807.86
                                                                                                                   112542.58
                                                                                                                                  0|
        3 15619304
                         Oniol
                                      502
                                             France|Female| 42|
                                                                    8 | 159660.8
                                                                                                           0|
                                                                                                                   113931.57
                                                                                                                                  1|
        41
           15701354
                         Boni l
                                      6991
                                             France|Female| 39|
                                                                       9.91
                                                                                                                    93826.63
                                                                                                                                  0|
        5| 15737888|Mitchell|
                                      8501
                                              Spain|Female| 43|
                                                                    2 1 1 2 5 5 1 0 . 8 2 1
                                                                                                                     79084.1
                                                                                                                                  0
```

Full outer join

```
|RowNumber|CustomerId| | Surname|CreditScore|Geography|Gender|Age|Tenure| | Balance|NumOfProducts|IsActiveMember|EstimatedSalary|Exited|
     1288 | 15565701 | Ferri
                                      698
                                              Spain|Female| 39|
                                                                    9|161993.89|
                                                                                                                    90212.38
                                                                                                           0|
                                                                                                                                  0
           15565706 | Akobundu |
                                              Spain| Male| 35|
                                                                    1 0.0
     4199
                                      612
                                                                                                           1|
                                                                                                                    83256.26
                                                                                                                                  1|
                                            France | Male | 47
                                                                    1 64430.06
     7091
            15565714 | Cattaneo |
                                      601
                                                                                                                    96517.97
                                                                                                                                  0|
                                      627 | Germany | Female | 30 |
     2021
           15565779 Kent
                                                                    6| 57809.32|
                                                                                                                    188258.49
                                                                                                                                   0|
                                                                                                            0|
                                      745 | Germany | Male | 48 |
      3698 | 15565796 | Docherty |
                                                                   10 | 96048.55 |
                                                                                                                    74510.65
                                                                                                                                   0|
```

```
# 3.Left join
spark.sql("""

SELECT cc.*
FROM credit_card cc
left JOIN new n
ON n.CustomerId = cc.CustomerId
""").show(5)
```

```
|RowNumber| CustomerId | Surname | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | IsActiveMember | EstimatedSalary | Exited | Surname | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | IsActiveMember | EstimatedSalary | Exited | Geography | Gender | Geography | G
                                                                                                                                                                                                    619| France|Female| 42|
                                             1 15634602 Hargrave
                                                                                                                                                                                                                                                                                                                                                                                                0.01
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     11
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         101348.88
                                                                                                                                                                                                                                           Spain|Female| 41|
                                                              15647311
                                                                                                                                 H6111
                                                                                                                                                                                                    6081
                                                                                                                                                                                                                                                                                                                                                          1 83807.86
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    11
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         112542.58
                                                                                                                                                                                                                                      France|Female| 42|
France|Female| 39|
                                                                15619304
                                                                                                                                 Oniol
                                                                                                                                                                                                    502 l
                                                                                                                                                                                                                                                                                                                                                          8 | 159660.8
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     3|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                          113931.57
                                             41
                                                              15701354
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            93826.631
                                                                                                                                 Bonil
                                                                                                                                                                                                    699 l
                                                                                                                                                                                                                                                                                                                                                                                          0.01
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    21
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    0|
                                             5| 15737888|Mitchell|
                                                                                                                                                                                                                                            Spain|Female| 43|
                                                                                                                                                                                                                                                                                                                                                           2|125510.82|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 79084.1
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                    øl
                                                                                                                                                                                                    850 l
only showing top 5 rows
```

Right join

```
# Right join
spark.sql("""
    SELECT cc.*
    FROM credit_card cc
    Right JOIN new n
    ON n.CustomerId = cc.CustomerId
""").show(5)
```

```
|RowNumber| CustomerId | Surname | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | IsActiveMember | EstimatedSalary | Exited | Surname | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | IsActiveMember | EstimatedSalary | Exited | Geography | Gender | Geography | G
                                           1 | 15634602 | Hargrave |
                                                                                                                                                                                             619| France|Female| 42|
                                                                                                                                                                                                                                                                                                                                                                                  0.0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       101348.88
                                                             15647311
                                                                                                                             Hill|
                                                                                                                                                                                              608
                                                                                                                                                                                                                                    Spain|Female| 41|
                                                                                                                                                                                                                                                                                                                                               1 83807.86
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       112542.58
                                                            15619304
                                                                                                                             Onio|
                                                                                                                                                                                             502
                                                                                                                                                                                                                               France|Female| 42|
                                                                                                                                                                                                                                                                                                                                               8| 159660.8|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       113931.57
                                           41
                                                             15701354
                                                                                                                             Boni l
                                                                                                                                                                                              699 l
                                                                                                                                                                                                                                France|Female| 39|
                                                                                                                                                                                                                                                                                                                                                                              0.01
                                                                                                                                                                                                                                                                                                                                                                                                                                                                     21
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           93826.63
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                0|
                                           5 | 15737888 | Mitchell |
                                                                                                                                                                                              8501
                                                                                                                                                                                                                                  Spain|Female| 43|
                                                                                                                                                                                                                                                                                                                                               2|125510.82|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               79084.1
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               0
only showing top 5 rows
```

Left Anti join

```
#Left Anti join
spark.sql("""

SELECT cc.*
FROM credit_card cc
LEFT ANTI JOIN new n
ON n.CustomerId = cc.CustomerId
""").show(5)
```

```
|RowNumber| CustomerId| Surname | CreditScore| Geography | Gender| Age | Tenure| Balance | NumOfProducts | IsActive Member| Estimated Salary | Exited| Salary | Salary | Exited| Salary | Sala
                                                                15634602 | Hargrave |
                                                                                                                                                                                                                                                 France|Female| 42|
                                                                                                                                                                                                                                                                                                                                                                                                              0.0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   101348.88
                                                                15647311
                                                                                                                                                                                                                                                                                                                                                                          1 83807.86
                                                                                                                                                                                                          608
                                                                                                                                                                                                                                                  Spain|Female| 41|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   112542.58
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  0|
                                                                15619304
                                                                                                                                      Onio
                                                                                                                                                                                                          502
                                                                                                                                                                                                                                                France|Female| 42|
                                                                                                                                                                                                                                                                                                                                                                          8 | 159660.8
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                   113931.57
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                        0
                                                                                                                                                                                                                                                 France|Female| 39|
                                                                15701354
                                                                                                                                      Bonil
                                                                                                                                                                                                          699 l
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       93826.631
                                                                                                                                                                                                                                                                                                                                                                                                           0.01
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                         0
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  0|
                                            5 15737888 Mitchell
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  0
                                                                                                                                                                                                                                                                                                                                                                          2|125510.82|
                                                                                                                                                                                                          850 l
                                                                                                                                                                                                                                                      Spain|Female| 43|
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             79084.1
```

Left semi join

```
#left semi join
spark.sql("""

    SELECT cc.*
    FROM credit_card cc
    LEFT SEMI JOIN new n
    ON n.CustomerId = cc.CustomerId
""").show(5)
```

```
|RowNumber|CustomerId| Surname|CreditScore|Geography|Gender|Age|Tenure| Balance|NumOfProducts|IsActiveMember|EstimatedSalary|Exited|
        6 15574012 Chu
                                   645| Spain| Male| 44|
                                                                    8 | 113755.78 |
                                                                                                                                  1|
                                                                                                          0|
                                  822| France| Male| 50|
376| Germany|Female| 29|
        7| 15592531|Bartlett|
                                     822| France| Male| 50|
                                                                                                                    10062.8
                                                                           0.0
                                                                                                                                  0|
           15656148 Obinna
                                                                    4 | 115046.74 |
                                                                                                                   119346.88
                                   501| France| Male| 44|
684| France| Male| 27|
           15792365
                           He|
                                                                    4|142051.07|
                                                                                                                    74940.5
                                                                                                                                  0|
       10 | 15592389
                                                                    2|134603.88|
                                                                                                                    71725.73
```

3. Simple Aggregations

Calculate total Balance Amount of all customers

Calculate revenue of Balance Amount(Calculate average)

4.Group By

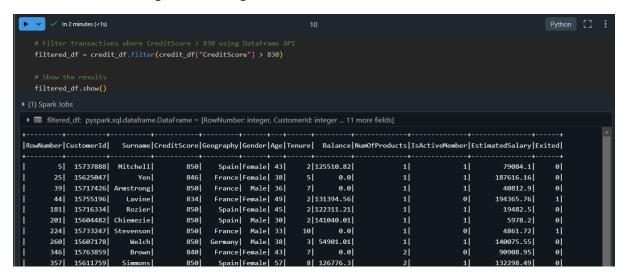
Calculate revenue genereated in each geographical region

Using PySpark

Initializing Pyspark session

1.Filter

List all customers having credit score greater tan 830



List all female customers

```
► ■ female_customers_df: pyspark.sql.dataframe.DataFrame = [RowNumber: integer, Customerld: integer ... 11 more fields]
|RowNumber|CustomerId| Surname|CreditScore|Geography|Gender|Age|Tenure| Balance|NumOfProducts|IsActiveMember|EstimatedSalary|Exited|
        1 | 15634602 | Hargrave
                                            France|Female| 42|
                                                                    1 83807.86
           15647311
                                              Spain|Female| 41|
                                                                                                                  112542.58
           15619304
                          Onio|
                                       502
                                             France|Female| 42|
                                                                    8 159660.8
                                                                                                                  113931.57
           15701354
                          Bonil
                                       699 l
                                             France|Female| 39|
                                                                           0.01
                                                                                                          0|
                                                                                                                  93826.63
                                                                                                                                0|
                                                                    2 | 125510.82 |
           15737888 | Mitchell|
                                       850 l
                                              Spain|Female| 43|
                                                                                                          1
                                                                                                                   79084.1
                                            Germany|Female| 29|
                                                                                                                  119346.88
       8
           15656148
                        Obinnal
                                       376 l
                                                                    4|115046.74|
                                                                                                          0|
       13
           15632264
                          Kay
                                             France Female 34
                                                                   10
                                                                            0.0
                                                                                                          0
                                                                                                                  26260.98
                                                                                                                                0|
                                       476
       14
           15691483
                          Chin|
                                       549
                                             France|Female| 25|
                                                                            0.0
                                                                                                          0|
                                                                                                                  190857.79
                                                                                                                                0
           15600882
                         Scott|
                                       635
                                              Spain|Female| 35|
                                                                                                                  65951.65
                                                                                                                                0|
                                                                            0.0
                                              Spain|Female| 24|
                                                                                                                   14406.41
                                       549
                                                                            0.0
       20|
            15568982
                           Hao|
                                       726
                                              France|Female| 24|
                                                                            0.0
                                                                                                                   54724.03
           15597945| Dellucci|
                                       636
                                              Spain|Female| 32|
                                                                            0.0
                                                                                                                  138555.46
       23
           15699309| Gerasimov|
                                       510
                                              Spain|Female| 38|
                                                                            0.0
                                                                                                          0|
                                                                                                                  118913.53
       25
           15625047
                           Yenl
                                       846
                                             France|Female| 38|
                                                                            0.01
                                                                                                          1
                                                                                                                  187616.16
                                                                                                                                0|
            15728693 McWilliams
```

2.Joins

```
#Joins

joined_df = credit_df.join(new_df, credit_df["CustomerId"] == new_df["CustomerId"], "inner").show()

joined_df = credit_df.join(new_df, credit_df["CustomerId"] == new_df["CustomerId"], "outer").show()

joined_df = credit_df.join(new_df, credit_df["CustomerId"] == new_df["CustomerId"], "left").show()

joined_df = credit_df.join(new_df, credit_df["CustomerId"] == new_df["CustomerId"], "right").show()

joined_df = credit_df.join(new_df, credit_df["CustomerId"] == new_df["CustomerId"], "leftanti").show()

joined_df = credit_df.join(new_df, credit_df["CustomerId"] == new_df["CustomerId"], "leftsemi").show()
```

Innerjoin output

```
|RowNumber|CustomerId| Surname|CreditScore|Geography|Gender|Age|Tenure| Balance|NumOfProducts|IsActiveMember|EstimatedSalary|Exited|RowNumber|C
ustomerId| Surname|CreditScore|Geography|Gender|Age|Tenure| Balance|NumOfProducts|IsActiveMember|EstimatedSalary|Exited|
       1 | 15634602 | Hargrave |
                                  619 France Female | 42 | 2 |
                                                                     คดไ
                                                                                                         101348.88
15634602 | Hargrave |
                      619| France|Female| 42| 2| 0.0|
                                                                                             101348.88
      2 15647311
                                  608 | Spain|Female| 41|
                                                              1 83807.86
                                                                                                         112542.58
                                                                                                                      0 l
                      Hill|
15647311| Hill|
                              Spain|Female| 41| 1| 83807.86|
                                                                                             112542.58
                                                                                                         0
                       608
                                                                      1
       3 15619304
                                   502| France|Female| 42|
                                                             8 159660.8
                                                                                                         113931.57
                                                                                                0|
                       Onio
15619304 Onio
                             France|Female| 42| 8| 159660.8|
                                                                                             113931.57
       4 15701354
                                   699| France|Female| 39|
                                                                                                          93826.63
                       Boni|
15701354 Boni
                       699 l
                             France|Female| 39| 1| 0.0|
                                                                                              93826.63
                              850| Spain|Female| 43| 2|125510.82|
Spain|Female| 43| 2|125510.82| 1
       5| 15737888|Mitchell|
                                                                                                           79084.1
.
15737888|Mitchell|
                                                                                               79084.1
                       850
```

Full outer join output

```
|RowNumber|CustomerId| Surname|CreditScore|Geography|Gender|Age|Tenure| Balance|NumOfProducts|IsActiveMember|EstimatedSalary|Exited|RowNumber
[CustomerId] Surname | CreditScore | Geography | Gender | Age | Tenure | Balance | NumOfProducts | IsActiveMember | EstimatedSalary | Exited | Salary | Exited | Ex
              1288 15565701
                                                                    Ferri|
                                                                                                        698
                                                                                                                        Spain|Female| 39|
                                                                                                                                                                                      9 | 161993.89 |
                                                                                                                                                                                                                                                                                                                    90212.38
                                                                                                                                                                                                                                                                                                                                                       0|
                                                                                                                                                                                                                                                                                                                                                                         nul1
               null
                                 null|
                                                                    null|
                                                                                               null| null|null| null|
                                                                                                                                                                  null|
                                                                                                                                                                                                         null|
                                                                                                                                                                                                                                                null|
                                                                                                                                                                                                                                                                                         null|
                                                                                                                                                                                                                                                                                                           null|
                               15565706
                                                                                                      612 | Spain | Male | 35 |
                                                                                                                                                                                      11
                                                                                                                                                                                                           0.01
                                                                                                                                                                                                                                                                                          1|
                                                                                                                                                                                                                                                                                                                  83256.26
                                                                                                                                                                                                                                                                                                                                                                         nul1
              4199
                                                            Akobundu l
                                                                   null|
                                                                                               null| null|null| null|
                                                                                                                                                                                                                                                                                                           null|
                                  null|
                                                                                                                                                                 null|
                                                                                                                                                                                                         null|
                                                                                                                                                                                                                                                null|
                                                                                                                                                                                                                                                                                         null|
               null|
                                                                                                                                                                                                                                                                                                                   96517.97
                                                                                                                                                                                                                                                                                                                                                                         nu11
                               15565714
                                                                                                        601| France| Male| 47|
                                                                                                                                                                                      1 64430.06
                                                                                                                                                                                                                                                                                                                                                       0|
              7091
                                                            Cattaneo
                                                                    null|
                                                                                               null| null|null| null|
               null|
                                  null|
                                                                                                                                                                  null|
                                                                                                                                                                                                        null|
                                                                                                                                                                                                                                                null|
                                                                                                                                                                                                                                                                                         null|
                                                                                                         627 | Germany | Female | 30 |
                                                                                                                                                                                                                                                                                                                188258.49
              2021
                                                                                               null| null|null| null|
                                                                                                                                                                                                                                                 null|
               3698|
                               15565796
                                                                                                         745| Germany| Male| 48|
                                                                                                                                                                                                                                                                                                                  74510.65
                                                                                                                                                                                                                                                                                                                                                                          null
                 null|
                                                                    null|
                                                                                               null| null|null| null|
                                                                                                                                                                   null|
                                                                                                                                                                                                          null|
                                                                                                                                                                                                                                                 null|
                                                                                                                                                                                                                                                                                          null|
                                                                                                                                                                                                                                                                                                           null|
                                                                                                                                                                                                                                                                                                                  30583.95
               3417 l
                               155658061
                                                                  Toosey
                                                                                                        532 France Male 38
                                                                                                                                                                                                           0.01
                                                                                                                                                                                                                                                  2|
                                                                                                                                                                                                                                                                                          0|
                                                                                                                                                                                                                                                                                                                                                       0|
                                                                                                                                                                                                                                                                                                                                                                         null
                                                                                               null| null|null| null|
                                  null|
                                                                                                                                                                  null|
                                                                                                                                                                                                                                                                                                           null|
               null
                                                                    null|
                                                                                                                                                                                                         null|
                                                                                                                                                                                                                                                null|
                                                                                                                                                                                                                                                                                         null|
                               15565878
                                                                    Bates
                                                                                                      631 | Spain | Male | 29
                                                                                                                                                                                       3|
                                                                                                                                                                                                           0.0
                                                                                                                                                                                                                                                                                                                197963.46
                                                                                                                                                                                                                                                                                                                                                       0|
                                                                                                                                                                                                                                                                                                                                                                         null
              6882
                                                                                               null| null|null|
                                     null|
                                                                     null|
                                                                                                                                               null|
                                                                                                                                                                                                                                                                                          null|
                                                                                                                                                                                                                                                                                                             null|
```

Left join output

 	+-		+-				+		+-
		reditScore Geography Gender Age Tenu						ited RowN	umber C
		Geography Gender Age Tenure Balance							
1 15634602 Ha	rgrave	619 France Female 42	2	0.0	1	1	101348.88	1	1
15634602 Hargrave	619	France Female 42 2 0.0		1	1	101348.88	1		
2 15647311	Hill	608 Spain Female 41	1	83807.86	1	1	112542.58	0	2
15647311 Hill	608	Spain Female 41 1 83807.86		1	1	112542.58	0		
3 15619304	Onio	502 France Female 42	8	159660.8	3	0	113931.57	1	3
15619304 Onio	502	France Female 42 8 159660.8		3	0	113931.57	1		
4 15701354	Boni	699 France Female 39	1	0.0	2	0	93826.63	0	4
15701354 Boni	699	France Female 39 1 0.0		2	0	93826.63	0		
5 15737888 Mi	tchell	850 Spain Female 43	2 1	25510.82	1	1	79084.1	0	5
15737888 Mitchell	850	Spain Female 43 2 125510.82		1	1	79084.1	0		
++	+-		+-		+		+	+	

Right join output

```
e|CreditScore|Geography|Gender|Age|Tenure| Balance|NumOfProducts|IsActiveMember|EstimatedSalary|Exited|RowN
ustomerId| Surname|CreditScore|Geography|Gender|Age|Tenure| Balance|NumOfProducts|IsActiveMember|EstimatedSalary|Exited|
        1 15634602 Hargrave
                                     619 France Female | 42
                                                                                                                101348.88
                                                                  2|
                                                                          0.01
15634602 | Hargrave |
                               France|Female| 42|
                                                                                                   101348.88
                        619
                                     608 | Spain|Female| 41|
       2 15647311
                                                                                                                112542.58
15647311
                         608|
                                 Spain|Female| 41|
                                                     1 83807.86
                                                                                                   112542.58
            15619304
                        Onio|
                                                                                                                113931.57
15619304
            Onio|
                        502
                               France|Female| 42|
                                                     8 | 159660.8
                                                                                           0|
                                                                                                   113931.57
        4 | 15701354
                        Boni|
                                     699| France|Female| 39|
                                                                  1|
                                                                          0.0
                                                                                                                 93826.63
                                                                                                                                        4|
15701354
            Boni|
                        699 l
                                France|Female| 39|
                                                            0.0
                                                                                           0|
                                                                                                    93826.63
                                     850| Spain|Female| 43|
        5 | 15737888 | Mitchell |
                                                                  2|125510.82|
                                                                                                                  79084.1
                                                                                                                              0
                                Spain|Female| 43|
                                                                                                     79084.1
15737888|Mitchell|
                                                     2 | 125510.82 |
                                                                                                                 0|
                        850
```

Left anti join output

F	RowNumber	 CustomerId	Surname	CreditScore 0	Geography	Gender	 Age	Tenure	Bala	nce	NumOfProducts	IsActiveMember	 EstimatedSalary	Exited
+-	+			+-			+			+			++	
ш	6	15574012	Chu	645	Spain	Male	44	8	113755	5.78	2	0	149756.71	1
1	7	15592531	Bartlett	822	France	Male	50	7		0.0	2	1	10062.8	0
П	8	15656148	Obinna	376	Germany	Female	29	4	115046	5.74	4	0	119346.88	1
Ī.	9	15792365	He	501	France	Male	44	4	142051	1.07	2	1	74940.5	0
П	10	15592389	н?	684	France	Male	27	2	13460	8.88	1	1	71725.73	0
П	11	15767821	Bearce	528	France	Male	31	6	102016	5.72	2	0	80181.12	0
П	12	15737173	Andrews	497	Spain	Male	24	3		0.0	2	0	76390.01	0
П	13	15632264	Kay	476	France	Female	34	10		0.0	2	0	26260.98	0
Ī.	14	15691483	Chin	549	France	Female	25	5		0.0	2	0	190857.79	0
П	15	15600882	Scott	635	Spain	Female	35	7		0.0	2	1	65951.65	0
1	16	15643966	Goforth	616	Germany	Male	45	3	143129	.41	2	1	64327.26	0
П	17	15737452	Romeo	653	Germany	Male	58	1	132602	2.88	1	0	5097.67	1
1	181	157882181	Handancanl	5/10	Snainl	Eomalo	امدا	اه		ام م	2	1	14406 41	اه

Left semi join output

```
|RowNumber|CustomerId| Surname|CreditScore|Geography|Gender|Age|Tenure| Balance|NumOfProducts|IsActiveMember|EstimatedSalary|Exited|
             15634602 | Hargrave |
                                               France|Female| 42|
                                                                               0.0
             15647311
                          Hill|
                                                 Spain|Female| 41|
                                                                        1| 83807.86|
                                                                                                                           112542.58
             15619304
                          Onio|
                                        502
                                                France|Female| 42|
                                                                        8 | 159660.8
                                                                                                                           113931.57
                                               France|Female| 39|
Spain|Female| 43|
                                                                                                                  0|
1|
         4|
5|
             15701354
                          Boni
                                         699
                                                                                0.0
                                                                                                  2|
1|
                                                                                                                            93826.63
                                                                                                                                          0|
0|
             15737888 Mitchell
                                                                        2 | 125510.82 |
                                        850 l
                                                                                                                            79084.1
```

3. Simple Aggregations

Calucate sum and avg of balance amount of all customers

```
from pyspark.sql.functions import sum, avg
# 1. Total Balance Amount
total_balance = credit_df.agg(sum("Balance").alias("TotalBalance"))
total_balance.show()

# 2. Average of Balance Amount
average_balance = credit_df.agg(avg("Balance").alias("AverageBalance"))
average_balance.show()

* (4) Spark Jobs

* *** total_balance: pyspark.sql.dataframe.DataFrame = [TotalBalance: double]

* *** average_balance: pyspark.sql.dataframe.DataFrame = [AverageBalance: double]

* *** TotalBalance|

* *** TotalBalance|

* *** TotalBalance|

* *** AverageBalance|

* *** TotalBalance|

* TotalBal
```

4.Group by

Calculate average transaction amount of each location