



Batch Data Processing with Apache Spark

Tolgahan Cepel - Mert Akat

Contents

1. [Introduction](#)
2. [Importing the libraries](#)
3. [Reading the data](#)
4. [SparkSQL API Practices](#)
 - [Selecting columns](#)
 - [Data manipulation](#)
 - [Filtering rows](#)
 - [Aggregating data](#)
 - [Joining](#)
5. [Case Studies](#)
 - [Assignment 1: Jacket sales per region](#)
 - [Assignment 2: Maximum turnover of the retailer regions](#)

1. Introduction



SQL Tables Description

- **FactSale:** Sales transactions fact table
- **FactPurchase:** Purchases fact table
- **DimRetailer:** Retailer details dimension table
- **DimCustomer:** Customer details dimension table
- **DimProduct:** Product details dimension table
- **DimRegion:** Region details dimension table
- **DimDate:** Date dimension table
- **DimSupplier:** Supplier details dimension table

2. Importing the libraries

```
In [1]: 1 from pyspark.sql import SparkSession, functions as F
```

3. Reading the data

```
In [51]: 1 # Creating new SparkSession instance
2 spark = SparkSession.builder.appName("PySparkExample").getOrCreate()
3
4 # Reading parquet data and assigning to DataFrame variables
5 df_pur = spark.read.parquet("data/purchase")
6 df_sal = spark.read.parquet("data/sale")
7 df_cus = spark.read.parquet("data/customer")
8 df_ret = spark.read.parquet("data/retailer")
9 df_pro = spark.read.parquet("data/product")
10 df_sup = spark.read.parquet("data/supplier")
11 df_reg = spark.read.parquet("data/region")
12 df_date = spark.read.parquet("data/date")
13
14 # Creating temporary view tables for Spark SQL queries
15 df_cus.createOrReplaceTempView("DimCustomer")
16 df_pur.createOrReplaceTempView("FactPurchase")
17 df_sal.createOrReplaceTempView("FactSale")
18 df_ret.createOrReplaceTempView("DimRetailer")
19 df_pro.createOrReplaceTempView("DimProduct")
20 df_sup.createOrReplaceTempView("DimSupplier")
21 df_reg.createOrReplaceTempView("DimRegion")
22 df_date.createOrReplaceTempView("DimDate")
```

4. Spark SQL Practices

Selecting columns

```
In [3]: 1 spark.sql("SELECT customer_id, name, surname, birth_date FROM DimCustom
```

```
+-----+-----+-----+-----+
|customer_id|  name| surname|birth_date|
+-----+-----+-----+-----+
|          1| Jazmin|  Burril|1958-09-22|
|          2| Dalila|   Faers|2000-11-08|
|          3|Wayland|Walework|1976-03-08|
|          4|Amberly|  Haquin|1948-10-08|
|          5|Garrett|  Frear|1957-09-25|
+-----+-----+-----+-----+
```

In [4]: 1 df_cus.select("customer_id", "name", "surname", "birth_date").show(5)

```
+-----+-----+-----+-----+
|customer_id|  name| surname|birth_date|
+-----+-----+-----+-----+
|          1| Jazmin|  Burril|1958-09-22|
|          2| Dalila|   Faers|2000-11-08|
|          3|Wayland|Walework|1976-03-08|
|          4|Amberly|  Haquin|1948-10-08|
|          5|Garrett|  Frear|1957-09-25|
+-----+-----+-----+-----+
only showing top 5 rows
```

Data manipulation: Calculating the ages from date of birth data.

In [5]: 1 spark.sql("""
2 SELECT
3 customer_id
4 ,name
5 ,surname
6 ,YEAR(CURRENT_DATE()) - YEAR(birth_date) AS age
7 FROM DimCustomer
8 LIMIT 5
9 """).show()

```
+-----+-----+-----+-----+
|customer_id|  name| surname|age|
+-----+-----+-----+-----+
|          1| Jazmin|  Burril| 66|
|          2| Dalila|   Faers| 24|
|          3|Wayland|Walework| 48|
|          4|Amberly|  Haquin| 76|
|          5|Garrett|  Frear| 67|
+-----+-----+-----+-----+
```

In [6]: 1 (
2 df_cus.withColumn("age", F.year(F.current_date()) - F.year("birth_date"))
3 .select("customer_id", "name", "surname", "age")
4 .show(5)
5)

```
+-----+-----+-----+-----+
|customer_id|  name| surname|age|
+-----+-----+-----+-----+
|          1| Jazmin|  Burril| 66|
|          2| Dalila|   Faers| 24|
|          3|Wayland|Walework| 48|
|          4|Amberly|  Haquin| 76|
|          5|Garrett|  Frear| 67|
+-----+-----+-----+-----+
only showing top 5 rows
```

Filtering rows

In [7]:

```

1 spark.sql("""
2 SELECT
3     name
4     ,surname
5     ,age
6 FROM
7 (
8     SELECT
9         customer_id
10        ,name
11        ,surname
12        ,YEAR(CURRENT_DATE()) - YEAR(birth_date) AS age
13    FROM DimCustomer
14 )
15 WHERE age >= 30
16 LIMIT 5
17 """).show()

```

```

+-----+-----+-----+
|  name| surname|age|
+-----+-----+-----+
| Jazmin|  Burril| 66|
|Wayland|Walework| 48|
|Amberly|  Haquin| 76|
|Garrett|  Frear| 67|
|  Horst|  Isted| 49|
+-----+-----+-----+

```

In [8]:

```

1 (
2     df_cus.withColumn("age", F.year(F.current_date()) - F.year("birth_c
3     .select("name", "surname", "age")
4     .filter(F.col("age") >= 30)
5     .show(5)
6 )

```

```

+-----+-----+-----+
|  name| surname|age|
+-----+-----+-----+
| Jazmin|  Burril| 66|
|Wayland|Walework| 48|
|Amberly|  Haquin| 76|
|Garrett|  Frear| 67|
|  Horst|  Isted| 49|
+-----+-----+-----+
only showing top 5 rows

```

Aggregating data

In [9]:

```

1 spark.sql("""
2 SELECT
3     order_id
4     ,SUM(quantity) AS total_quantity
5     ,SUM(total_amt) AS total_amount
6 FROM FactSale
7 GROUP BY order_id
8 ORDER BY total_quantity DESC
9 LIMIT 10
10 """).show()

```

```

+-----+-----+-----+
|order_id|total_quantity|total_amount|
+-----+-----+-----+
|    3647|           13|         521|
|    2574|           13|         488|
|    3515|           13|         402|
|     101|           12|         359|
|     440|           12|         426|
|    3763|           12|         323|
|    1585|           12|         488|
|    3289|           12|         327|
|    1382|           11|         452|
|    1752|           11|         298|
+-----+-----+-----+

```

In [10]:

```

1 (
2     df_sal.groupBy("order_id").agg(
3         F.sum("quantity").alias("total_quantity"),
4         F.sum("total_amt").alias("total_amount")
5     ).orderBy("total_quantity", ascending=False)
6     .show(10)
7 )

```

```

+-----+-----+-----+
|order_id|total_quantity|total_amount|
+-----+-----+-----+
|    3647|           13|         521|
|    2574|           13|         488|
|    3515|           13|         402|
|     101|           12|         359|
|     440|           12|         426|
|    3763|           12|         323|
|    1585|           12|         488|
|    3289|           12|         327|
|    2337|           11|         357|
|    3743|           11|         359|
+-----+-----+-----+

```

only showing top 10 rows

Joining

In [11]:

```

1 spark.sql("""
2 SELECT
3     region_name
4     ,AVG(YEAR(CURRENT_DATE()) - YEAR(birth_date)) AS age
5 FROM DimCustomer cus
6 INNER JOIN DimRegion reg
7 ON cus.city_id = reg.city_id
8 GROUP BY region_name
9 ORDER BY age DESC
10 """).show()

```

region_name	age
Akdeniz	50.81521739130435
Dogu Anadolu	50.13095238095238
Guneydogu Anadolu	48.58119658119658
Marmara	48.189542483660134
Ic Anadolu	48.07772020725388
Karadeniz	47.75121951219512
Ege	46.888888888888886

In [12]:

```

1 (
2     df_cus
3     .join(df_reg, df_cus.city_id == df_reg.city_id)
4     .groupBy("region_name").agg(
5         F.avg(F.year(F.current_date()) - F.year("birth_date")).alias("a
6     )
7     .orderBy("age", ascending=False)
8     .show()
9 )

```

region_name	age
Akdeniz	50.81521739130435
Dogu Anadolu	50.13095238095238
Guneydogu Anadolu	48.58119658119658
Marmara	48.189542483660134
Ic Anadolu	48.07772020725388
Karadeniz	47.75121951219512
Ege	46.888888888888886

5. Case Studies

Assignment 1: Jacket sales per region

Write SparkSQL and Python API scripts that results: Region-based total quantity and amount of jacket sales between June and August 2023.

The expected out is as follows:

region_name	product_type	total_quantity	total_amount
Marmara	Jacket	213	8358
Dogu Anadolu	Jacket	284	11547
Guneydogu Anadolu	Jacket	176	6981
Ic Anadolu	Jacket	260	10496
Akdeniz	Jacket	162	6637
Karadeniz	Jacket	310	12582
Ege	Jacket	101	3953

External links

- <https://spark.apache.org/docs/3.1.2/api/python/reference/api/pyspark.sql.DataFrame.join>.
(<https://spark.apache.org/docs/3.1.2/api/python/reference/api/pyspark.sql.DataFrame.join>)
- <https://spark.apache.org/docs/latest/api/python/reference/pyspark.sql/api/pyspark.sql.DataFrame.groupBy>.
(<https://spark.apache.org/docs/latest/api/python/reference/pyspark.sql/api/pyspark.sql.DataFrame.groupBy>)
- <https://spark.apache.org/docs/3.1.1/api/python/reference/api/pyspark.sql.DataFrame.groupBy>.
(<https://spark.apache.org/docs/3.1.1/api/python/reference/api/pyspark.sql.DataFrame.groupBy>)



In [25]:

```

1  # Your Spark SQL Solution:
2  from pyspark.sql import SparkSession
3  from pyspark.sql.functions import month, year, sum, lit
4  # Filtering jacket sales between June and August 2023
5  jacket_sales = spark.sql("""
6      SELECT
7          reg.region_name,
8          prod.product_type,
9          SUM(sal.quantity) AS total_quantity,
10         SUM(sal.total_amt) AS total_amount
11     FROM
12         FactSale sal
13     INNER JOIN
14         DimProduct prod ON sal.product_id = prod.product_id
15     INNER JOIN
16         DimCustomer cus ON sal.customer_id = cus.customer_id
17     INNER JOIN
18         DimRegion reg ON cus.city_id = reg.city_id
19     INNER JOIN
20         DimDate dat ON sal.date = dat.date
21     WHERE
22         prod.product_type = 'Jacket'
23         AND dat.month BETWEEN 6 AND 8
24         AND dat.year = 2023
25     GROUP BY
26         reg.region_name, prod.product_type
27 """)
28
29
30 # Displaying the result
31 jacket_sales.show()
32

```

region_name	product_type	total_quantity	total_amount
Marmara	Jacket	213	8358
Dogu Anadolu	Jacket	284	11547
Guneydogu Anadolu	Jacket	176	6981
Ic Anadolu	Jacket	260	10496
Akdeniz	Jacket	162	6637
Karadeniz	Jacket	310	12582
Ege	Jacket	101	3953

In [31]:

```

1  # Your PySpark Solution:
2
3
4  from pyspark.sql import functions as F
5
6  # Ceket satışlarını filtreleme ve bölgeye göre gruplama
7  jacket_sales_per_region = df_sal.join(df_pro, df_sal.product_id == df_p
8      .join(df_cus, df_sal.customer_id == df_cus.customer_id) \
9      .join(df_reg, df_cus.city_id == df_reg.city_id) \
10     .join(df_date, df_sal.date == df_date.date) \
11     .filter((df_pro.product_type == 'Jacket') &
12             (F.month(df_date.date).between(6, 8)) &
13             (df_date.year == 2023)) \
14     .groupBy("region_name", "product_type") \
15     .agg(F.sum("quantity").alias("total_quantity"),
16         F.sum("total_amt").alias("total_amount"))
17
18 # Sonucu gösterme
19 jacket_sales_per_region.show()
20

```

```

+-----+-----+-----+-----+
| region_name | product_type | total_quantity | total_amount |
+-----+-----+-----+-----+
| Marmara | Jacket | 213 | 8358 |
| Dogu Anadolu | Jacket | 284 | 11547 |
| Guneydogu Anadolu | Jacket | 176 | 6981 |
| Ic Anadolu | Jacket | 260 | 10496 |
| Akdeniz | Jacket | 162 | 6637 |
| Karadeniz | Jacket | 310 | 12582 |
| Ege | Jacket | 101 | 3953 |
+-----+-----+-----+-----+

```

Assignment 2: Maximum turnover of the retailer regions

Find the maximum turnover region of each retailer, and obtain total amount for each retailer and region.

The expected out is as follows:

retailer_id	retailer_name	region_name	total_amount
1	Hepsiburada	Karadeniz	42642
2	Trendyol	Ic Anadolu	71689
3	n11	Ic Anadolu	11995
4	Gittigidiyor	Karadeniz	16081

In [62]:

```

1  # Your Spark SQL Solution:
2
3
4  max_turnover_per_retailer = spark.sql("""
5      SELECT
6          retailer_id,
7          retailer_name,
8          region_name,
9          total_amount
10     FROM (
11         SELECT
12             retailer_id,
13             retailer_name,
14             region_name,
15             total_amount,
16             ROW_NUMBER() OVER(PARTITION BY retailer_id ORDER BY total_a
17     FROM (
18         SELECT
19             ret.retailer_id,
20             ret.retailer_name,
21             reg.region_name,
22             SUM(sal.total_amt) AS total_amount
23     FROM
24         FactSale sal
25     JOIN
26         DimRetailer ret ON sal.retailer_id = ret.retailer_id
27     JOIN
28         DimCustomer cus ON sal.customer_id = cus.customer_id
29     JOIN
30         DimRegion reg ON cus.city_id = reg.city_id
31     GROUP BY
32         ret.retailer_id, ret.retailer_name, reg.region_name
33     ) turnover_per_region
34     ) ranked_turnover
35     WHERE rn = 1
36 """)
37
38 # Displaying the result
39 max_turnover_per_retailer.show()
40
41

```

```

+-----+-----+-----+-----+
|retailer_id|retailer_name|region_name|total_amount|
+-----+-----+-----+-----+
|          1| Hepsiburada|  Karadeniz|         42642|
|          2|   Trendyol|  Ic Anadolu|         71689|
|          3|         n11|  Ic Anadolu|         11995|
|          4| Gittigidiyor|  Karadeniz|         16081|
+-----+-----+-----+-----+

```

In [40]:

```

1  # Your PySpark Solution:
2  from pyspark.sql import functions as F
3  from pyspark.sql.window import Window
4
5  # Joining necessary tables with DataFrame aliases
6  joined_df = df_sal.alias("sal").join(df_ret.alias("ret"), F.col("sal.re
7  .join(df_cus.alias("cus"), F.col("sal.customer_id") =
8  .join(df_reg.alias("reg"), F.col("cus.city_id") == F.
9
10 # Grouping data by retailer_id, region_name, and aggregating total amou
11 grouped_df = joined_df.groupBy("ret.retailer_id", "ret.retailer_name",
12 .agg(F.sum("total_amt").alias("total_amount"))
13
14 # Defining a window partitioned by retailer_id and ordered by total_amo
15 window_spec = Window.partitionBy("ret.retailer_id").orderBy(F.desc("tot
16
17 # Adding a rank column to find the maximum turnover region for each ret
18 result_df = grouped_df.withColumn("rank", F.rank().over(window_spec)) \
19 .filter(F.col("rank") == 1) \
20 .drop("rank")
21
22 # Renaming columns and ordering the result
23 result_df = result_df.select("ret.retailer_id", "ret.retailer_name", "r
24 .orderBy("ret.retailer_id")
25
26 # Showing the result
27 result_df.show()
28
29

```

```

+-----+-----+-----+-----+
|retailer_id|retailer_name|region_name|total_amount|
+-----+-----+-----+-----+
|          1|  Hepsiburada|  Karadeniz|         42642|
|          2|    Trendyol|  Ic Anadolu|         71689|
|          3|          n11|  Ic Anadolu|         11995|
|          4| Gittigidiyor|  Karadeniz|         16081|
+-----+-----+-----+-----+

```