

GA-4

Report: Implementing SCD Type II on a Customer Master Data Frame using PySpark

Assignment Overview:

For this assignment, we are required to implement Slowly Changing Dimension (SCD) Type II on a customer master data frame using PySpark. The implementation will be executed on a Google Cloud Dataproc cluster. The input files will be created manually based on the example data discussed in the lecture.

Objective:

The objective is to manage historical changes in a customer master data frame by maintaining historical records with a start and end date for each record version.

Implementation Steps:

1. **Set Up the Environment:**
 - Set up a Google Cloud Dataproc cluster.
 - Install PySpark on the cluster.
2. **Create Input Data Files:**
 - Create two CSV files, `original.csv` and `updated.csv`, representing the customer master data and the updated records, respectively.
3. **Write PySpark Code:**
 - Initialize a Spark session.
 - Read the input data files into PySpark DataFrames.
 - Implement the logic to update existing records and add new records with proper version control.
 - Save the final DataFrame to an output file.

Input Data:

original.csv:

```
idx,name,dob,start_date,end_date
1,Mary,16-04-2001,01-01-1970,10-12-2077
2,Joseph,21-07-2002,01-01-1970,10-12-2077
3,Robert,06-11-2001,01-01-1970,10-12-2077
```

updated.csv:

```
name,dob
Mary,17-09-2002
Andrew,13-02-2004
```

PySpark Code:

```
from datetime import datetime
from pyspark.sql import SparkSession
from pyspark.sql.functions import when, lit, col, max as max_
from pyspark.sql.types import StructType, StructField, IntegerType,
StringType

# Current date
current_date = datetime.now().strftime("%d-%m-%Y")

# Create a Spark session
spark = SparkSession.builder.appName("SCD_Type_2").getOrCreate()

# Define schema for the input data
schema = StructType([
    StructField("idx", IntegerType(), True),
    StructField("name", StringType(), True),
    StructField("dob", StringType(), True),
    StructField("start_date", StringType(), True),
    StructField("end_date", StringType(), True)
])

# Access data
original = spark.read.csv("original.csv", header=True, schema=schema)
updated = spark.read.csv("updated.csv", header=True, inferSchema=True)

# Update end_date for matching records
for row in updated.collect():
    condition1 = col("name") == row.name
    condition2 = col("end_date") > current_date

    original = original.withColumn("end_date", when(condition1 &
condition2, lit(current_date)).otherwise(col("end_date")))
```

```

# Add new rows from updated DataFrame to original DataFrame
max_idx = original.agg(max_("idx")).collect()[0][0] if
original.count() > 0 else 0

new_rows = updated.withColumn("idx", lit(max_idx + 1)) \
                  .withColumn("start_date", lit(current_date)) \
                  .withColumn("end_date", lit("10-12-2077"))

original = original.union(new_rows.select(original.columns))

# Show the data
original.show()

# Save to CSV
output_csv_path = "output.csv"
original.toPandas().to_csv(output_csv_path, index=False, header=True)

# Save to text file
output_txt_path = "output.txt"
with open(output_txt_path, 'w') as f:
    original_data = original.collect()
    for row in original_data:

f.write(f"{row.idx}\t{row.name}\t{row.dob}\t{row.start_date}\t{row.end
_date}\n")

# Stop the Spark session
spark.stop()

```

Output Data:

The resulting DataFrame after applying SCD Type II is as follows:

idx	name	dob	start_date	end_date
1	Mary	16-04-2001	01-01-1970	10-12-2077
2	Joseph	21-07-2002	01-01-1970	10-12-2077

3	Robert	06-11-2001	01-01-1970	10-12-2077
4	Mary	17-09-2002	14-07-2024	10-12-2077
5	Andrew	13-02-2004	14-07-2024	10-12-2077

Execution Instructions:

1. Upload the `original.csv` and `updated.csv` files to the Google Cloud Storage bucket.
2. Submit the PySpark code to the Google Cloud Dataproc cluster for execution.
3. Download the output files `output.csv` and `output.txt` from the Cloud Storage bucket.

Submission:

Submit the final output files along with the PySpark code and this report through the provided

Conclusion:

This implementation of SCD Type II in PySpark maintains historical records and enables tracking of changes over time in the customer master data frame. The solution ensures that each change is captured with a new record version, providing a complete history of customer data changes.