

Προχωρημένα Θέματα Βάσεων Δεδομένων

Εργαστηριακή Αναφορά

Ευθύμιος Καραγιάννης (03119434)

Χειμώνιο 2023-2024

Generate Dataframe

```
1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import udf, col
3 from pyspark.sql.types import IntegerType, DoubleType, DateType
4 from datetime import date
5
6 APP_NAME = "GenerateDataFrame"
7 HDFS_DATA_DIR = "hdfs://oceanos-master:54310/data"
8
9 spark = SparkSession.builder.appName(APP_NAME).getOrCreate()
10
11
12 def convert_date_format(input_date: str):
13     date_part, -, - = input_date.split("_")
14     month, day, year = date_part.split("/")
15     return date(int(year), int(month), int(day))
16
17
18 convert_date_udf = udf(convert_date_format, DateType())
19
20 csv_file_path1 = f"{HDFS_DATA_DIR}/Crime-Data-from-2010-to-2019.csv"
21 csv_file_path2 = f"{HDFS_DATA_DIR}/Crime-Data-from-2020-to-Present.csv"
22
23 df1 = spark.read.csv(csv_file_path1, header=True)
24 df2 = spark.read.csv(csv_file_path2, header=True)
25
26 df = df1.union(df2)
```

```

27
28
29 df = (
30     df.withColumn("Date_Rptd", convert_date_udf(col("Date_Rptd")))
31     .withColumn("DATE_OCC", convert_date_udf(col("DATE_OCC")))
32     .withColumn("Vict_Age", df["Vict_Age"].cast(IntegerType()))
33     .withColumn("LAT", df["LAT"].cast(DoubleType()))
34     .withColumn("LON", df["LON"].cast(DoubleType()))
35 )
36
37 output_csv_path = f"{HDFS_DATA_DIR}/Full_Data"
38 df.write.option("header", True).mode("overwrite").csv(output_csv_path)

```

Query 1

```

1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import month, year, row_number, desc
3 from pyspark.sql.window import Window
4 import time
5
6 APP_NAME = "DF_Query_1"
7 HDFS_DATA_DIR = "hdfs://okeanos-master:54310/data"
8 HDFS_OUTPUT_DIR = "hdfs://okeanos-master:54310/results"
9
10 spark = SparkSession.builder.appName(APP_NAME).getOrCreate()
11
12 csv_file_path = f"{HDFS_DATA_DIR}/Full_Data"
13
14 df = spark.read.csv(csv_file_path, header=True)
15
16 start_time = time.time()
17
18 window = Window.partitionBy("year").orderBy(desc("crime_total"))
19 results_df = (
20     df.withColumn("year", year("Date_Rptd"))
21     .withColumn("month", month("Date_Rptd"))
22     .groupBy("year", "month")
23     .count()
24     .withColumnRenamed("count", "crime_total")
25     .withColumn(
26         "rank",
27         row_number().over(window),
28     )
29     .filter("rank <= 3")

```

```

30         .select("year", "month", "crime_total", "rank")
31     )
32
33     results_df.show(results_df.count(), truncate=False)
34     exec_time = time.time() - start_time
35     print(f"\n\nExec_time:_{exec_time}_sec")
36
37     spark.stop()

```

```

1  from pyspark.sql import SparkSession
2  from time import time
3
4  APP_NAME = "SQL_Query_1"
5  HDFS_DATA_DIR = "hdfs://okeanos-master:54310/data"
6
7  spark = SparkSession.builder.appName(APP_NAME).getOrCreate()
8
9  csv_file_path = f"{HDFS_DATA_DIR}/Full_Data"
10 df = spark.read.csv(csv_file_path, header=True)
11
12 start_time = time()
13
14 df.createOrReplaceTempView("tmp_view")
15 sql_query = """
16     SELECT year, month, crime_total, rank
17     FROM (
18         SELECT
19             year('Date Rptd') AS year,
20             month('Date Rptd') AS month,
21             COUNT(*) AS crime_total,
22             ROWNUMBER() OVER (PARTITION BY year('Date Rptd') ORDER BY
23                 COUNT(*) DESC) AS rank
24     FROM
25         tmp_view
26     GROUP BY
27         year('Date Rptd'), month('Date Rptd')
28     ) tmp
29     WHERE rank <= 3
30 """
31 results_df = spark.sql(sql_query)
32
33 results_df.show(results_df.count(), truncate=False)
34 exec_time = time() - start_time
35 print(f"\n\nExec_time:_{exec_time}_sec\n\n")

```

Εκτελούμε τα παραπάνω queries στον master ή στον worker node με τις παρακάτω εντολές:

spark-submit --num-executors 4 ./src/query_1.py

spark-submit --num-executors 4 ./src/sql_query_1.py

Παρακάτω παρουσιάζονται τα αποτελέσματα των παραπάνω queries καθώς και χρόνοι εκτέλεσης των SQL και DataFrame API. Όσον αφορά τους χρόνους εκτέλεσης, και τα δύο APIs εμφανίζουν σχεδόν τα ίδια αποτελέσματα, αφού και τα δύο χρησιμοποιούν το Spark SQL engine και τον Catalyst optimizer. Οι παρακάτω χρόνοι αποτελούν τον μέσο χρόνο 5 εκτελέσεων.

Table 1: Query 1 results

Year	Month	Crime Total	Rank
2010	3	12828	1
2010	7	12085	2
2010	4	12070	3
2011	8	20643	1
2011	5	20643	2
2011	10	20533	3
2012	10	30963	1
2012	8	30646	2
2012	5	29842	3
2013	1	11429	1
2013	3	8229	2
2013	8	8056	3
2014	1	6297	1
2014	6	5712	2
2014	5	5694	3
2015	7	10106	1
2015	5	10102	2
2015	3	10065	3
2016	10	16405	1
2016	12	15796	2
2016	8	15761	3
2017	3	27409	1
2017	8	27216	2
2017	7	26858	3
2018	1	8357	1
2018	2	6167	2
2018	8	6135	3
2019	7	19127	1
2019	8	18874	2
2019	5	18538	3

2020	1	7130	1
2020	2	5818	2
2020	7	5210	3
2021	7	28426	1
2021	10	27955	2
2021	8	27690	3
2022	8	32162	1
2022	7	31372	2
2022	10	30936	3
2023	8	26308	1
2023	10	26249	2
2023	7	25974	3
2024	1	1	1

Table 2: Query 1 times

API	Time (sec)
SQL	28.55
DataFrame	32.04

Query 2

```

1 from pyspark.sql import SparkSession
2 from pyspark.sql.functions import col, when, desc
3 from time import time
4
5 APP_NAME = "DF_Query_2"
6 HDFS_DATA_DIR = "hdfs://oceanos-master:54310/data"
7
8 spark = SparkSession.builder.appName(APP_NAME).getOrCreate()
9
10 csv_file_path = f"{HDFS_DATA_DIR}/Full_Data"
11
12 df = spark.read.csv(csv_file_path, header=True)
13
14 start_time = time()
15
16 results_df = (
17     df.filter(col("Premis_Desc") == "STREET")

```

```

18     .withColumn(
19         "day_split",
20         when((col("TIME_OCC") >= "0500") & (col("TIME_OCC") < "1200"),
21             "morning")
22         .when((col("TIME_OCC") >= "1200") & (col("TIME_OCC") <
23             "1700"), "afternoon")
24         .when((col("TIME_OCC") >= "1700") & (col("TIME_OCC") <
25             "2100"), "evening")
26         .otherwise("night"),
27     )
28     .groupBy("day_split")
29     .count()
30     .withColumnRenamed("count", "crime_total")
31     .orderBy(desc("crime_total"))
32 )
33 results_df.show(results_df.count(), truncate=False)
34 exec_time = time() - start_time
35 print(f"\n\nExec_time: {exec_time}_sec\n\n")
36 spark.stop()

```

```

1 from pyspark import SparkContext
2 from pyspark.sql import SparkSession
3 from time import time
4
5 APP_NAME = "RDD_Query_2"
6 HDFS_DATA_DIR = "hdfs://oceanos-master:54310/data"
7
8 sc = SparkContext(appName=APP_NAME)
9 spark = SparkSession(sc)
10
11 csv_file_path = f"{HDFS_DATA_DIR}/Full_Data"
12 init_rdd = spark.read.csv(csv_file_path, header=True).rdd
13
14
15 def map_time_to_day_split(row):
16     time_occ = row[3]
17     if "0500" <= time_occ < "1200":
18         return "morning", 1
19     elif "1200" <= time_occ < "1700":
20         return "afternoon", 1
21     elif "1700" <= time_occ < "2100":
22         return "evening", 1
23     else:

```

```

24         return "night", 1
25
26
27 start_time = time()
28
29 res_rdd = (
30     init_rdd.filter(lambda row: row[15] == "STREET")
31     .map(map_time_to_day_split)
32     .reduceByKey(lambda x, y: x + y)
33     .sortBy(lambda x: x[1], ascending=False)
34 )
35
36 exec_time = time() - start_time
37 print(f"\n\nExec_time: {exec_time}_sec\n\n")
38
39 results = res_rdd.collect()
40 for row in results:
41     print(row[0], row[1])
42
43 sc.stop()

```

Υποβάλλουμε τις παραπάνω εργασίες στο spark με τις παρακάτω εντολές:

spark-submit --num-executors 4 ./src/query_2.py

spark-submit --num-executors 4 ./src/rdd_query_2.py

Στην συνέχεια παρουσιάζονται τα αποτελέσματα και οι χρόνοι εκτέλεσης των DataFrame και RDD API. Στην περίπτωση αυτή βλέπουμε τα δύο APIs να παρουσιάζουν αρκετά διαφορετικούς χρόνους εκτέλεσης, με το RDD να είναι σημαντικά πιο αργό. Αυτό συμβαίνει διότι τα RDDs αποτελούν περισσότερο μία περιγραφή της λύσης και δεν μπορούν να γίνουν optimized από το Spark, καθώς δεν είναι γνωστός ο τύπος των δεδομένων και οι πράξεις μεταξύ τους. Επιπλέον απελούν αντικείμενα στην μνήμη του JVM, με αποτελέσματα να εξαρτώνται σε μεγάλο βαθμό από αυτό και να προστίθεται επιπλέον overhead όταν ο όγκος των δεδομένων μεγαλώνει (Garbage Collection, Object Serialization).

Table 3: Query 2 results

Day Split	Crime Total
Night	223171
Evening	174170
Afternoon	139129
Morning	117035

Table 4: Query 2 times

API	Time (sec)
RDD	91.02
DataFrame	27.00

Query 3

```

1 from pyspark.sql import SparkSession
2 from pyspark.sql.types import StringType
3 from pyspark.sql.functions import col, year, desc, udf, split,
  regexp_replace
4 from time import time
5
6 APP_NAME = "DF_Query_3"
7 HDFS_DATA_DIR = "hdfs://okeanos-master:54310/data"
8
9 spark = SparkSession.builder.appName(APP_NAME).getOrCreate()
10
11 basic_csv_path = f"{HDFS_DATA_DIR}/Full_Data"
12 revgecoding_csv_path = f"{HDFS_DATA_DIR}/revgecoding.csv"
13 income_csv_path = f"{HDFS_DATA_DIR}/LA_income_2015.csv"
14
15 basic_df = spark.read.csv(basic_csv_path, header=True)
16 revgecoding_df = spark.read.csv(revgecoding_csv_path, header=True)
17 income_df = spark.read.csv(income_csv_path, header=True)
18
19 victim_descent_mapping = {
20     "A": "Other_Asian",
21     "B": "Black",
22     "C": "Chinese",
23     "D": "Cambodian",
24     "F": "Filipino",
25     "G": "Guamanian",
26     "H": "Hispanic/Latin/Mexican",
27     "I": "American_Indian/Alaskan_Native",
28     "J": "Japanese",
29     "K": "Korean",
30     "L": "Laotian",
31     "O": "Other",
32     "P": "Pacific_Islander",
33     "S": "Samoan",

```



```

34     "U": "Hawaiian",
35     "V": "Vietnamese",
36     "W": "White",
37     "X": "Unknown",
38     "Z": "Asian_Indian",
39 }
40
41 mapping_udf = udf(lambda x: victim_descent_mapping.get(x, x),
42                   StringType())
43 results = []
44 start_time = time()
45
46 inner_join_df = (
47     basic_df.filter((year(col("Date_Rptd")) == 2015) & (col("Vict_
48         Descent") != ""))
49     .join(
50         revgecoding_df,
51         (basic_df["LAT"] == revgecoding_df["LAT"])
52         & (basic_df["LON"] == revgecoding_df["LON"]),
53     )
54     .select("Vict_Descent", "ZIPcode")
55 )
56 income_df = (
57     income_df.withColumn(
58         "Estimated_Median_Income", regexp_replace("Estimated_Median_
59             Income", "\\$", ""))
60     .withColumn(
61         "Estimated_Median_Income", regexp_replace("Estimated_Median_
62             Income", ",", ""))
63     .withColumn(
64         "Estimated_Median_Income", col("Estimated_Median_
65             Income").cast("integer"))
66 )
67
68 for asc in [True, False]:
69     sorted_income_ZIP_codes_df = (
70         income_df.sort("Estimated_Median_Income", ascending=asc)
71         .limit(3)
72         .select("Zip_Code")
73         .collect()

```

```

74     )
75
76     data = [row["Zip_Code"] for row in sorted_income_ZIP_codes_df]
77
78     results_df = (
79         inner_join_df.filter(split(col("ZIPcode"),
80             "_").getItem(0).isin(data))
81         .withColumnRenamed("Vict_Descent", "Victim_Descent")
82         .groupBy("Victim_Descent")
83         .count()
84         .withColumnRenamed("count", "crime_total")
85         .orderBy(desc("crime_total"))
86         .withColumn("Victim_Descent", mapping_udf(col("Victim_Descent"))))
87
88     results.append(results_df)
89
90
91 for df in results:
92     df.show(df.count(), truncate=False)
93     print()
94 exec_time = time() - start_time
95 print(f'\n\nExec_time: {exec_time}_sec\n\n')
96
97 spark.stop()

```

Για το ερώτημα αυτό, τα queries για τις περιοχές με το υψηλότερο και χαμηλότερο οικεγενειακό εισόδημα εκτελέστηκαν χωριστά, καθώς με αυτό το τρόπο τα αποτελέσματα παρουσιάζουν περισσότερο ενδιαφέρον. Αυτήν την φορά υποβάλλουμε την εργασία, για διαφορετικό πλήθος spark-executors, με την παρακάτω εντολή:

spark-submit --num-executors n ./src/query_3.py όπου **n = 2, 3, 4**.

Αυξάνοντας τον αριθμό των spark-executors, αυξάνουμε τον αριθμό των tasks της υποβολής που μπορούν να εκτελεστούν παράλληλα. Επιπλέον, κάθε executor δεσμεύει πόρους μνήμης και CPU τους κόμβου, ώστε να αποθηκεύει τοπικά τα δεδομένα και συνεπώς να εκτελεί τις πράξεις γρηγορότερα. Επομένως, αύξηση του πλήθους τους μπορεί να οδηγήσει σε καλύτερες επιδόσεις. Ωστόσο υπάρχουν περιπτώσεις όπου τα δεδομένα δεν είναι κατανομημένα ομοιόμορφα στους κόμβους και η αύξηση της παραλληλίας δεν οφελεί, ή τα tasks είναι πολύ μικρά και η δρομολόγηση τους σε executors προσθέτει σημαντικό overhead σε σχέση με την επεξεργασία των δεδομένων, με αποτελέσματα την αύξηση του χρόνου εκτέλεσης.

Table 5: Query 3 results (Table A)

Victim Descent	Crime Total
Hispanic/Latin/Mexican	144
White	116
Black	101
Unknown	39
Other	32
Other Asian	6
Chinese	1

Table 6: Query 3 results (Table B)

Victim Descent	Crime Total
Other	2
Other Asian	1
White	1
Hispanic/Latin/Mexican	1

Table 7: Query 3 times

Spark-executors	Time (sec)
2	48.46
3	43.95
4	53.98

Query 4

```

1 from pyspark.sql import SparkSession , Window
2 from pyspark.sql.types import FloatType
3 from pyspark.sql.functions import (
4     col ,
5     length ,
6     udf ,
7     year ,
8     avg ,
9     count ,
10    desc ,

```

```

11     min as spark_min ,
12 )
13 from geopy.distance import geodesic
14
15 APP_NAME = "DF_Query_4"
16 HDFS_DATA_DIR = "hdfs://okeanos-master:54310/data"
17
18 spark = SparkSession.builder.appName(APP_NAME).getOrCreate()
19
20 basic_csv_path = f"{HDFS_DATA_DIR}/Full_Data"
21 departmens_csv_path = f"{HDFS_DATA_DIR}/LAPD_Police_Stations.csv"
22
23 basic_df = spark.read.csv(basic_csv_path, header=True)
24 departments_df = spark.read.csv(departmens_csv_path, header=True)
25
26
27 def compute_distance(lat1, lon1, lat2, lon2):
28     crime_location = (lat1, lon1)
29     department_location = (lat2, lon2)
30     return geodesic(crime_location, department_location).kilometers
31
32
33 compute_distance_udf = udf(compute_distance, FloatType())
34
35 results_df_list = []
36
37 filtered_df = basic_df.filter(
38     col("Weapon_Used_Cd").startswith("1")
39     & (length(col("Weapon_Used_Cd")) == 3)
40     & ((col("LAT") > 0) & (col("LON") < 0))
41 )
42
43 filtered_df = filtered_df.withColumn("AREA",
44     filtered_df["AREA"].cast("integer"))
45 departments_df = departments_df.withColumn(
46     "PREC", departments_df["PREC"].cast("integer")
47 )
48 inner_join_df = filtered_df.join(departments_df, col("AREA") ==
49     col("PREC")).withColumn(
50     "distance",
51     compute_distance_udf(
52         filtered_df["LAT"],
53         filtered_df["LON"],
54         departments_df["Y"],

```

```

54         departments_df["X"],
55     ),
56 )
57
58 window_spec = Window.partitionBy("DR_NO").orderBy("distance")
59
60 cross_join_df = (
61     filtered_df.crossJoin(
62         departments_df.withColumnRenamed("LOCATION", "department_
63             location")
64     )
65     .withColumn(
66         "distance",
67         compute_distance_udf(
68             filtered_df["LAT"],
69             filtered_df["LON"],
70             departments_df["Y"],
71             departments_df["X"],
72         ),
73     )
74     .withColumn("min_distance",
75         spark_min("distance").over(window_spec))
76     .filter(col("distance") == col("min_distance"))
77     .drop("min_distance")
78 )
79
80 results_df_list = []
81
82 for join_df in [inner_join_df, cross_join_df]:
83     results_df = (
84         join_df.withColumn("year", year("Date_Rptd"))
85         .groupBy("year")
86         .agg(
87             avg("distance").alias("average_distance"),
88             count("*").alias("total_crimes")
89         )
90         .select("year", "average_distance", "total_crimes")
91         .orderBy("year")
92     )
93     results_df_list.append(results_df)
94
95     results_df = (
96         join_df.withColumnRenamed("DIVISION", "division")
97         .groupBy("division")
98         .agg(

```

```

96         avg("distance").alias("average_distance"),
97         count("*").alias("total_crimes")
98     )
99     .select("division", "average_distance", "total_crimes")
100     .orderBy(desc("total_crimes"))
101 )
102 results_df_list.append(results_df)
103 for df in results_df_list:
104     df.show(df.count(), truncate=False)
105     print()
106
107 spark.stop()

```

Για να υποβάλλουμε την εργασία 4 στο spark χρειάζεται αρχικά να δημιουργήσουμε ένα virtual environment και να εγκαταστήσουμε τη βιβλιοθήκη geopy. Στη συνέχεια παράγουμε το αρχείο requirements.txt, το οποίο θα χρησιμοποιηθεί για την δημιουργία του dependencies.zip.

```

python3 -m venv venv
source venv/bin/activate
pip3 install geopy
pip3 freeze > requirements.txt
pip3 install -t dependencies -r requirements.txt
zip dependencies.zip dependencies/*

```

Τέλος, υποβάλλουμε την εργασία στο spark με την εντολή:

```
spark-submit --py-files dependencies.zip ./src/query_4.py
```

Table 8: Query 4 (1a)

Year	Average Distance	Total Crimes
2010	2.667	5842
2011	2.835	8186
2012	2.834	11956
2013	2.747	2322
2014	2.681	2354
2015	2.655	3505
2016	2.681	6496
2017	2.762	9883
2018	2.577	2282
2019	2.741	7100
2020	2.489	2287
2021	2.652	14415
2022	2.627	15722
2023	2.624	11767

Table 9: Query 4 (1b)

Division	Average Distance	Total Crimes
77TH STREET	2.657	16372
SOUTHEAST	2.103	10100
NEWTON	2.018	9469
SOUTHWEST	2.694	8867
HOLLENBECK	2.674	6382
HARBOR	4.065	5709
RAMPART	1.577	4994
NORTHEAST	3.860	4157
HOLLYWOOD	1.438	3946
OLYMPIC	1.830	3640
WILSHIRE	2.374	3639
MISSION	4.705	3592
CENTRAL	1.135	3470
WEST VALLEY	3.501	3270
FOOTHILL	3.824	3141
NORTH HOLLYWOOD	2.723	2968
VAN NUYS	2.214	2719
DEVONSHIRE	3.968	2209
PACIFIC	3.727	2193
TOPANGA	3.454	1666
WEST LOS ANGELES	4.232	1614

Table 10: Query 4 (2a)

Year	Average Distance	Total Crimes
2010	2.279	5842
2011	2.510	8186
2012	2.498	11956
2013	2.429	2322
2014	2.162	2354
2015	2.460	3505
2016	2.368	6496
2017	2.420	9883
2018	2.358	2282
2019	2.431	7100
2020	2.266	2287
2021	2.361	14415
2022	2.295	15722

2023 2.345 11764

Table 11: Query 4 (2b)

Division	Average Distance	Total Crimes
77TH STREET	1.688	12961
SOUTHWEST	2.282	11175
SOUTHEAST	2.237	9479
NEWTON	1.575	7044
WILSHIRE	2.456	6464
HOLLENBECK	2.665	6452
HOLLYWOOD	1.966	5793
HARBOR	3.871	5574
OLYMPIC	1.664	4814
RAMPART	1.409	4741
VAN NUYS	2.935	4580
FOOTHILL	3.596	3767
CENTRAL	1.011	3594
NORTHEAST	3.699	3337
NORTH HOLLYWOOD	2.763	2983
WEST VALLEY	2.726	2963
MISSION	3.819	2257
PACIFIC	3.706	2078
TOPANGA	3.031	1829
DEVONSHIRE	3.002	1164
WEST LOS ANGELES	2.686	1065

Join Strategies

```

1 from pyspark.sql import SparkSession
2 from pyspark.sql.types import StringType
3 from pyspark.sql.functions import col, year, desc, udf, split,
   regex_replace
4 from time import time
5 import sys
6
7 APP_NAME = "Join_Query_3"
8 HDFS_DATA_DIR = "hdfs://okeanos-master:54310/data"
9
10 spark = SparkSession.builder.appName(APP_NAME).getOrCreate()
11

```



```

12 basic_csv_path = f"{HDFS_DATA_DIR}/Full_Data"
13 revgecoding_csv_path = f"{HDFS_DATA_DIR}/revgecoding.csv"
14 income_csv_path = f"{HDFS_DATA_DIR}/LA_income_2015.csv"
15
16 basic_df = spark.read.csv(basic_csv_path, header=True)
17 revgecoding_df = spark.read.csv(revgecoding_csv_path, header=True)
18 income_df = spark.read.csv(income_csv_path, header=True)
19
20 victim_descent_mapping = {
21     "A": "Other_Asian",
22     "B": "Black",
23     "C": "Chinese",
24     "D": "Cambodian",
25     "F": "Filipino",
26     "G": "Guamanian",
27     "H": "Hispanic/Latin/Mexican",
28     "I": "American_Indian/Alaskan_Native",
29     "J": "Japanese",
30     "K": "Korean",
31     "L": "Laotian",
32     "O": "Other",
33     "P": "Pacific_Islander",
34     "S": "Samoan",
35     "U": "Hawaiian",
36     "V": "Vietnamese",
37     "W": "White",
38     "X": "Unknown",
39     "Z": "Asian_Indian",
40 }
41
42 mapping_udf = udf(lambda x: victim_descent_mapping.get(x, x),
43     StringType())
44 results = []
45
46 start_time = time()
47
48 revgecoding_df = revgecoding_df.hint(sys.argv[1])
49
50 inner_join_df = (
51     basic_df.filter((year(col("Date_Rptd")) == 2015) & (col("Vict_
52         Descent") != ""))
53     .join(
54         revgecoding_df,
55         (basic_df["LAT"] == revgecoding_df["LAT"])
56         & (basic_df["LON"] == revgecoding_df["LON"]),

```

```

55     )
56     .select("Vict_Descent", "ZIPcode")
57 )
58
59 income_df = (
60     income_df.withColumn(
61         "Estimated_Median_Income", regexp_replace("Estimated_Median_
62             Income", "\\$", "")
63     )
64     .withColumn(
65         "Estimated_Median_Income", regexp_replace("Estimated_Median_
66             Income", ",", "")
67     )
68     .withColumn(
69         "Estimated_Median_Income", col("Estimated_Median_
70             Income").cast("integer")
71     )
72 )
73
74 for asc in [True, False]:
75     sorted_income_ZIP_codes_df = (
76         income_df.sort("Estimated_Median_Income", ascending=asc)
77         .limit(3)
78         .select("Zip_Code")
79         .collect()
80     )
81
82 data = [row["Zip_Code"] for row in sorted_income_ZIP_codes_df]
83
84 results_df = (
85     inner_join_df.filter(split(col("ZIPcode"),
86         "_").getItem(0).isin(data))
87     .withColumnRenamed("Vict_Descent", "Victim_Descent")
88     .groupBy("Victim_Descent")
89     .count()
90     .withColumnRenamed("count", "crime_total")
91     .orderBy(desc("crime_total"))
92     .withColumn("Victim_Descent", mapping_udf(col("Victim_
93         Descent"))))
94 )
95
96 results.append(results_df)
97
98 for df in results:

```

```

95     df.explain()
96     df.show(df.count(), truncate=False)
97
98     exec_time = time() - start_time
99     print(f'\n\nExec_time:{exec_time}_sec\n\n')
100
101     spark.stop()

```

```

1  from pyspark.sql import SparkSession, Window
2  from pyspark.sql.types import FloatType
3  from pyspark.sql.functions import (
4      col,
5      length,
6      udf,
7      year,
8      avg,
9      count,
10     desc,
11 )
12 from geopy.distance import geodesic
13 import sys, time
14
15 APP_NAME = "Inner_Join_Query_4"
16 HDFS_DATA_DIR = "hdfs://oceanos-master:54310/data"
17
18 spark = SparkSession.builder.appName(APP_NAME).getOrCreate()
19
20 basic_csv_path = f"{HDFS_DATA_DIR}/Full_Data"
21 departmens_csv_path = f"{HDFS_DATA_DIR}/LAPD_Police_Stations.csv"
22
23 basic_df = spark.read.csv(basic_csv_path, header=True)
24 departments_df = spark.read.csv(departmens_csv_path, header=True)
25
26
27 def compute_distance(lat1, lon1, lat2, lon2):
28     crime_location = (lat1, lon1)
29     department_location = (lat2, lon2)
30     return geodesic(crime_location, department_location).kilometers
31
32
33 compute_distance_udf = udf(compute_distance, FloatType())
34
35 results_df_list = []
36
37 start_time = time.time()

```

```

38
39 filtered_df = basic_df.filter(
40     col("Weapon_Used_Cd").startswith("1")
41     & (length(col("Weapon_Used_Cd")) == 3)
42     & ((col("LAT") > 0) & (col("LON") < 0))
43 )
44
45 filtered_df = filtered_df.withColumn("AREA",
46     filtered_df["AREA"].cast("integer"))
47 departments_df = departments_df.withColumn(
48     "PREC", departments_df["PREC"].cast("integer")
49 )
50 departments_df = departments_df.hint(sys.argv[1])
51
52 join_df = filtered_df.join(departments_df, col("AREA") ==
53     col("PREC")).withColumn(
54     "distance",
55     compute_distance_udf(
56         filtered_df["LAT"],
57         filtered_df["LON"],
58         departments_df["Y"],
59         departments_df["X"],
60     ),
61 )
62 results_df_list = []
63
64 results_df = (
65     join_df.withColumn("year", year("Date_Rptd"))
66     .groupBy("year")
67     .agg(
68         avg("distance").alias("average_distance"),
69         count("*").alias("total_crimes")
70     )
71     .select("year", "average_distance", "total_crimes")
72     .orderBy("year")
73 )
74 results_df_list.append(results_df)
75
76 results_df = (
77     join_df.withColumnRenamed("DIVISION", "division")
78     .groupBy("division")
79     .agg(
80         avg("distance").alias("average_distance"),

```

```

80         count("*").alias("total_crimes")
81     ).select("division", "average_distance", "total_crimes")
82     .orderBy(desc("total_crimes"))
83 )
84 results_df_list.append(results_df)
85
86 for df in results_df_list:
87     df.explain()
88     df.show(df.count(), truncate=False)
89
90 exec_time = time.time() - start_time
91 print(f'\n\nExec_time:_{exec_time}_sec\n\n')
92
93 spark.stop()

```

Υποβάλλουμε τις εργασίες στο spark με τις παρακάτω εντολές:

spark-submit ./src/joins/query_3.py hint

spark-submit -py-files dependencies.zip ./src/joins/query_4.py hint

όπου **hint = broadcast, merge, shuffle_hash, shuffle_replicate_nl**.

Στις παραπάνω εργασίες μετρήθηκαν οι χρόνοι εκτέλεσης των ερωτημάτων 3 και 4 για διαφορετικές στρατηγικές join. Οι στρατηγικές που δοκιμάστηκαν είναι οι broadcast, merge, shuffle_hash και shuffle_replicate_nl. Παρατηρούμε ότι για το Query 3, οι στρατηγικές merge και shuffle_hash φαίνονται να παρουσιάζουν καλύτερα αποτελέσματα, ενώ στη περίπτωση του ερωτήματος 4, τα broadcast και shuffle_replicate_nl είναι αυτά με τους χαμηλότερους χρόνους. Αυτό οφείλεται στο γεγονός ότι στην δεύτερη περίπτωση, το departments dataset είναι αρκετά μικρό, επομένως το broadcast ή το replication του μικρότερου dataset σε κάθε executor έχει καλύτερη επίδοση από τα άλλα δύο joins, τα οποία εφαρμόζουν hashing ή sorting και shuffling στα δεδομένα. Αντίθετα, για την πρώτη περίπτωση, όπου το revgecoding dataset είναι σχετικά μεγαλύτερο, οι άλλες δύο στρατηγικές εμφανίζουν καλύτερα αποτελέσματα.

Table 12: Query 3 join times

Join	Time (sec)
Broadcast	59.06
Merge	47.88
Shuffle Hash	45.28
Shuffle Replicate NL	58.00

Table 13: Query 4 join times

Join	Time (sec)
Broadcast	61.87
Merge	83.36
Shuffle Hash	80.96
Shuffle Replicate NL	58.07