

Προχωρημένα Θέματα Βάσεων Δεδομένων

Εξαμηνιαία Εργασία

Κωνσταντίνος Πίκουλας (03120112) — konpik2002@gmail.com

Κωνσταντίνα Μαυρουκάκη-Καραγκούννη (03120063) – constantina.mavroukaki@gmail.com

Ιανουάριος 2025

Αποθετήριο [GitHub](#).

Περιεχόμενα

1	Ζητούμενο 1	2
2	Ζητούμενο 2	4
2.1	Υποερώτημα α	4
2.2	Υποερώτημα β	6
3	Ζητούμενο 3	8
4	Ζητούμενο 4	13
5	Ζητούμενο 5	16

1 Ζητούμενο 1

Στο ζητούμενο αυτό, καλούμαστε να χρησιμοποιήσουμε το DataFrame και το RDD API του PySpark, έτσι ώστε να ταξινομήσουμε σε φθίνουσα σειρά τα περιστατικά "Βαριάς Σωματικής Βλάβης", κατηγοριοποιημένα βάσει της ηλικιακής ομάδας στην οποία ανήκουν. Για την εκτέλεση του query, χρησιμοποιήθηκαν 4 Spark Executors.

Πρώτα απ' όλα, δημιουργήσαμε ένα SparkSession, και αντλήσαμε τα δεδομένα από τα δύο csv αρχεία που αφορούν τα εγκλήματα. Προκειμένου να είναι ενιαία προσπελάσιμα, εκτελέσαμε την πράξη union.

```
1 ## QUERY 1
2 APP_NAME = "Crime Victime Age Analysis"
3 SPARK_EXECUTORS = 4
4 spark = SparkSession.builder.appName(APP_NAME).config("spark.executor.instances", SPARK_EXECUTORS).
   getOrCreate()
5
6 # crime data
7 d1 = spark.read.csv(
8     "s3://initial-notebook-data-bucket-dblab-905418150721/CrimeData/
9     Crime_Data_from_2010_to_2019_20241101.csv",
10    header=True, inferSchema=True)
11 d2 = spark.read.csv(
12     "s3://initial-notebook-data-bucket-dblab-905418150721/CrimeData/
13     Crime_Data_from_2020_to_Present_20241101.csv",
14    header=True, inferSchema=True)
15 crime_data = d1.union(d2)
```

Αρχικά, χρησιμοποιούμε το DataFrame API προκειμένου να εκτελέσουμε το Query 1. Απομονώνουμε τα δεδομένα βάσει της περιγραφής του εγκλήματος, έτσι ώστε να κρατήσουμε μόνο δεδομένα που αφορούν "Βαριά Σωματική Βλάβη". Έπειτα, προσθέτουμε μια στήλη στα δεδομένα που αφορά το Age Group του θύματος, ως εξής:

- Παιδιά: < 18
- Νεαροί ενήλικοι: 18-24
- Ενήλικοι: 25-64
- Ηλικιωμένοι: > 64

Ομαδοποιούμε τα δεδομένα βάσει του Age Group και τέλος επιστρέφουμε τον αριθμό των θυμάτων για κάθε ηλικιακή ομάδα σε φθίνουσα σειρά.

```
1 # DATAFRAME API BEGIN #
2 start = time.time()
3 filtered_df = crime_data.filter(crime_data["Crm Cd Desc"].contains("AGGRAVATED ASSAULT")) #get only
   aggravated assault
4 grouped_df = filtered_df.withColumn("AGE GROUP", when(filtered_df['Vict Age'] <= 0, "Unknown").when(
   filtered_df['Vict Age'] < 18, "Children").when(filtered_df["Vict Age"] <= 24, "Young adults")
   .when(filtered_df['Vict Age'] <= 64, "Adults").when(filtered_df['Vict
   Age'] > 64, "Elderly"))
5
6 categories_df = grouped_df.groupBy("AGE GROUP").count().orderBy(desc('count')).collect() #group rows
   based on the age group assigned and count them
7 end = time.time()
8 dataframe_time = end-start
9 print("DATAFRAME performance:", end-start)
10 print("DATAFRAME results:", categories_df)
11 print("Dataframe Speedup:", rdd_time/dataframe_time)
12 # DATAFRAME API END #
```

Για το RDD API η διαδικασία είναι αρκετά παρόμοια. Δημιουργήσαμε μια συνάρτηση get_category_by_age, η οποία επιστρέφει την ηλικιακή ομάδα βάσει της ηλικίας. Αφού κρατήσαμε μόνο τα δεδομένα που χρειαζόμαστε, περάσαμε όλα τα δεδομένα από την συνάρτησή μας έτσι ώστε να κατηγοριοποιηθούν σωστά. Προκειμένου, όμως, να γίνει μετά η καταμέτρηση, φροντίσαμε τα δεδομένα να έχουν την μορφή (<Age Group>, 1), όπου <Age Group> η ηλικιακή ομάδα του θύματος βάσει της ηλικίας του. Έτσι, στο reduction βήμα, προσθέτουμε τα values (τους άσσους στην προκειμένη περίπτωση), παίρνοντας σαν αποτέλεσμα το πλήθος των θυμάτων ανά ηλικιακή ομάδα.

```
1 def get_category_by_age(age):
2     if age is None:
3         return "NULL"
4     age = int(age)
5     category = None
6     if age <= 0:
7         category = 'Unknown'
8     elif age < 18:
9         category = "Children"
```

```

10     elif age <= 24:
11         category = "Young adults"
12     elif age <= 64:
13         category = "Adults"
14     else:
15         category = "Elderly"
16     return category
17
18 # RDD API BEGIN #
19 start = time.time()
20 crime_rdd = crime_data.rdd # convert dataframe to rdd
21 filtered_rdd = crime_rdd.filter(lambda row: "AGGRAVATED ASSAULT" in row["Crm Cd Desc"])
22 grouped_rdd = filtered_rdd.map(lambda row: (get_category_by_age(row["Vict Age"]), 1))
23 categories_rdd = grouped_rdd.reduceByKey(lambda x,y: x+y).sortBy(lambda tup: -tup[1]).collect()
24 end = time.time()
25 rdd_time = end - start
26 print("RDD performance:", end - start)
27 print("RDD results:", categories_rdd)
28 # RDD API END #

```

Σημείωση: Η συνάρτηση `get_category_by_age` μπορούσε να χρησιμοποιηθεί και στο DataFrame API, ως User Defined Function. Εν τέλει, ο λόγος που δεν χρησιμοποιήθηκε, είναι επειδή πιθανότατα θα εμποδίσει το PySpark από το να εκτελέσει optimizations (περισσότερα μπορείτε να δείτε [εδώ](#)).

Παρατηρήσαμε πως το dataset περιέχει δεδομένα με θύματα που έχουν ηλικία ≤ 0 . Επιλέξαμε να δημιουργήσουμε μια επιπλέον κατηγορία `Unknown` για αυτά τα θύματα. Τα αποτελέσματα και των δύο APIs είναι προφανώς ίδια. Παρακάτω φαίνεται η έξοδος από το DataFrame API.

```

1 +-----+-----+
2 |   AGE GROUP | count |
3 +-----+-----+
4 |      Adults | 121093 |
5 | Young adults | 33605 |
6 |    Children | 10830 |
7 |    Elderly  | 5985 |
8 |    Unknown  | 5098 |
9 +-----+-----+

```

Το DataFrame API τελείωσε την εκτέλεση του σε 7.735968351364136 sec , ενώ το RDD API τελείωσε σε 16.28496217727661 sec. Παρατηρούμε δηλαδή ένα speedup ίσο με **2.1050967943017724** στην εκτέλεση του DataFrame API σε σχέση με το RDD API. Αυτή η διαφορά οφείλεται στο Catalyst, ένα optimization layer του Spark το οποίο χρησιμοποιείται στα DataFrames αλλά όχι στα RDDs. Έτσι, το Spark εκτελεί αρκετές βελτιστοποιήσεις κατά την εκτέλεση του query με το DataFrame API, με αποτέλεσμα την καλύτερη απόδοση του, σε σχέση με την εκτέλεση του query χρησιμοποιώντας το RDD API.

2 Ζητούμενο 2

2.1 Υποερώτημα α

Στο ερώτημα αυτό, καλούμαστε να βρούμε για κάθε έτος, τα τρία Αστυνομικά Τμήματα με το υψηλότερο ποσοστό κλεισμένων υποθέσεων και να παρουσιάσουμε τα αποτελέσματα με φθίνουσα σειρά ως προς το `closed_case_rate` και αύξουσα ως προς το έτος. Το query εκτελέστηκε τόσο με το DataFrame API, όσο και με το SQL API.

Τα δεδομένα αντλήθηκαν από τα csv με τα καταγεγραμμένα εγκλήματα, όπως προηγουμένως και χρησιμοποιήθηκε ξανά η πράξη `union`, έτσι ώστε να αντιμετωπίζεται ενιαία το dataset.

```
1 ### QUERY 2 (a)
2 APP_NAME = "Closed Cases"
3 spark = SparkSession.builder.appName(APP_NAME).getOrCreate() # reconstructing because we dont need 4
   spark executors now
4
5 # crime data
6 d1 = spark.read.csv(
7     "s3://initial-notebook-data-bucket-dblab-905418150721/CrimeData/
   Crime_Data_from_2010_to_2019_20241101.csv",
8     header=True, inferSchema=True).select(col('DATE OCC'), col('AREA NAME'), col('Status'), col('LAT
   '), col('LON'))
9 d2 = spark.read.csv(
10    "s3://initial-notebook-data-bucket-dblab-905418150721/CrimeData/
   Crime_Data_from_2020_to_Present_20241101.csv",
11    header=True, inferSchema=True).select(col('DATE OCC'), col('AREA NAME'), col('Status'), col('LAT
   '), col('LON'))
12 crime_data = d1.union(d2)
```

Παρακάτω φαίνεται η υλοποίηση του query χρησιμοποιώντας το DataFrame API. Γίνεται αρχικά εξαγωγή του έτους που καταγράφηκε το κάθε έγκλημα (σε νέα στήλη με το όνομα `YEAR`). Έπειτα, δημιουργείται μια νέα στήλη με το όνομα `closed` που συμβολίζει εάν η υπόθεση θεωρείται κλεισμένη ή όχι (χρησιμοποιώντας 1 ή 0 αντίστοιχα). Θεωρούμε μια υπόθεση κλεισμένη όταν η στήλη `Status` δεν έχει τις τιμές `IC` ή `UKN`. Στο επόμενο βήμα, κάνουμε `aggregate` τα rows σε ζευγάρια (`YEAR`, `AREA NAME`) και αθροίζουμε για το κάθε ζευγάρι την στήλη `closed` που δημιουργήσαμε προηγουμένως αλλά και το συνολικό πλήθος των rows. Έτσι, για κάθε Αστυνομικό Τμήμα και κάθε χρονιά έχουμε το πλήθος των κλεισμένων υποθέσεων αλλά και το πλήθος των υποθέσεων που ανέλαβε συνολικά. Υπολογίζουμε επομένως σε μια νέα στήλη `closed_case_rate`, το ποσοστό κλεισμένων υποθέσεων για κάθε ζευγάρι (`YEAR`, `AREA NAME`). Ορίζουμε, έπειτα, ένα `Window function` το οποίο κάνει `partition` το DataFrame μας βάσει του `YEAR`. Τα `partitioned` δεδομένα ταξινομούνται βάσει του ποσοστού κλεισμένων υποθέσεων και κρατούνται μόνο τα 3 υψηλότερα ποσοστά ανά χρονιά. Έτσι, για κάθε έτος επιστρέφουμε τα 3 υψηλότερα ποσοστά κλεισμένων υποθέσεων μαζί με τα αντίστοιχα Αστυνομικά Τμήματα.

```
1 # DATAFRAME API BEGIN #
2 start = time.time()
3 crime_data = crime_data.withColumn('YEAR', year(to_timestamp(col('DATE OCC'), 'MM/dd/yyyy hh:mm:ss a
   '))).withColumn('closed', when(~col('Status').isin('IC', 'UKN'), 1).otherwise(0))
4 crime_data_agg = crime_data.groupBy('YEAR', 'AREA NAME').agg(
5     count('*').alias("total_cases"),
6     sum('closed').alias('closed_cases')
7 )
8 crime_cc_rates = crime_data_agg.withColumn('closed_case_rate', when(col('total_cases') > 0, col('
   closed_cases')/col('total_cases') * 100).otherwise(None)).drop('closed_cases').drop('total_cases
   ')
9 window_spec = Window.partitionBy('year').orderBy(desc('closed_case_rate'))
10 ranked_df = crime_cc_rates.withColumn("#", row_number().over(window_spec))
11 top3_df = ranked_df.filter(col("#") <= 3)
12 top3_df_collect = top3_df.collect()
13 end = time.time()
14 print("DATAFRAME performance:", end-start)
15 print("DATAFRAME results:")
16 top3_df.show(n=top3_df.count(), truncate=False)
17 # DATAFRAME API END #
```

Για την εκτέλεση του query με SQL API, η διαδικασία είναι παρόμοια. Πρώτα γράφουμε το query σε SQL, όπου για την παραγωγή των ενδιάμεσων αποτελεσμάτων χρησιμοποιήθηκαν τα ανάλογα Common Table Expressions.

```
1 WITH processed_data AS (
2     SELECT
3         YEAR(TO_TIMESTAMP(`DATE OCC`, 'MM/dd/yyyy hh:mm:ss a')) AS YEAR,
4         `AREA NAME`,
5         CASE WHEN `Status` NOT IN ('IC', 'UKN') THEN 1 ELSE 0 END AS closed
6     FROM crime_data
7 ),
```

```

8 aggregated_data AS (
9     SELECT
10         YEAR,
11         `AREA NAME`,
12         COUNT(*) AS total_cases,
13         SUM(closed) AS closed_cases
14     FROM processed_data
15     GROUP BY YEAR, `AREA NAME`
16 ),
17 crime_cc_rates AS (
18     SELECT
19         YEAR,
20         `AREA NAME`,
21         (CASE WHEN total_cases > 0 THEN (closed_cases / total_cases) * 100 ELSE NULL END) AS
22         closed_case_rate
23     FROM aggregated_data
24 ),
25 ranked_data AS (
26     SELECT
27         YEAR,
28         `AREA NAME`,
29         closed_case_rate,
30         ROW_NUMBER() OVER (PARTITION BY YEAR ORDER BY closed_case_rate DESC) AS rank
31     FROM crime_cc_rates
32 )
33 SELECT
34     YEAR,
35     `AREA NAME`,
36     closed_case_rate,
37     rank AS `#`
38 FROM ranked_data
39 WHERE rank <= 3

```

Έπειτα εκτελούμε το SQL query μέσω του PySpark ως εξής:

```

1 # SQL API BEGIN #
2 start = time.time()
3 crime_data.createOrReplaceTempView("crime_data") # register the crime_data DataFrame as a temporary
4 view
5 top3_sql = spark.sql(query)
6 top3_sql_collect = top3_sql.collect()
7 end = time.time()
8 print("SQL performance:", end-start)
9 print("SQL results:")
10 top3_sql.show(n=top3_df.count(), truncate=False)
11 # SQL API END #

```

Τα αποτελέσματα και των δύο εκτελέσεων είναι προφανώς τα ίδια. Ενδεικτικά, παρακάτω φαίνεται η έξοδος από το DataFrame API.

```

1 +---+-----+-----+---+
2 |YEAR|AREA NAME |closed_case_rate |# |
3 +---+-----+-----+---+
4 |2010|Rampart    |32.84713448949121 |1 |
5 |2010|Olympic    |31.515289821999087|2 |
6 |2010|Harbor     |29.36028339237341 |3 |
7 |2011|Olympic    |35.040060090135206|1 |
8 |2011|Rampart    |32.4964471814306  |2 |
9 |2011|Harbor     |28.51336246316431 |3 |
10 |2012|Olympic    |34.29708533302119 |1 |
11 |2012|Rampart    |32.46000463714352 |2 |
12 |2012|Harbor     |29.509585848956675|3 |
13 |2013|Olympic    |33.58217940999398 |1 |
14 |2013|Rampart    |32.1060382916053  |2 |
15 |2013|Harbor     |29.723638951488557|3 |
16 |2014|Van Nuys   |32.0215235281705  |1 |
17 |2014|West Valley|31.49754809505847 |2 |
18 |2014|Mission    |31.224939855653567|3 |
19 |2015|Van Nuys   |32.265140677157845|1 |
20 |2015|Mission    |30.463762673676303|2 |
21 |2015|Foothill   |30.353001803658852|3 |
22 |2016|Van Nuys   |32.194518462124094|1 |
23 |2016|West Valley|31.40146437042384 |2 |
24 |2016|Foothill   |29.908647228131645|3 |
25 |2017|Van Nuys   |32.0554272517321  |1 |
26 |2017|Mission    |31.055387158996968|2 |

```

```

27 |2017|Foothill      |30.469700657094183|3 |
28 |2018|Foothill      |30.731346958877126|1 |
29 |2018|Mission       |30.727023319615913|2 |
30 |2018|Van Nuys      |28.905206942590123|3 |
31 |2019|Mission       |30.727411112319235|1 |
32 |2019|West Valley    |30.57974335472044 |2 |
33 |2019|N Hollywood    |29.23808669119627 |3 |
34 |2020|West Valley    |30.771131982204647|1 |
35 |2020|Mission       |30.14974649215894 |2 |
36 |2020|Harbor          |29.693486590038315|3 |
37 |2021|Mission       |30.318115590092276|1 |
38 |2021|West Valley    |28.971087440009363|2 |
39 |2021|Foothill      |27.993757094211126|3 |
40 |2022|West Valley    |26.536367172306498|1 |
41 |2022|Harbor          |26.337538060026098|2 |
42 |2022|Topanga        |26.234013317831096|3 |
43 |2023|Foothill      |26.76076020122974 |1 |
44 |2023|Topanga        |26.538022616453986|2 |
45 |2023|Mission       |25.662731120516817|3 |
46 |2024|N Hollywood    |19.598528961078763|1 |
47 |2024|Foothill      |18.620882188721385|2 |
48 |2024|77th Street    |17.586318167150694|3 |
49 +-----+-----+-----+-----+

```

Η έκδοση του query με το DataFrame API, εκτελείται σε 13.183066368103027 sec, ενώ η έκδοση με το SQL API, σε 8.88805389404297 sec. Παρατηρούμε, δηλαδή, πως η εκτέλεση της έκδοσης με το SQL API παρουσιάζει speedup **1.4831088982798082** σε σχέση με το DataFrame API. Ωστόσο χρονικά οι δύο εκδόσεις δεν απέχουν ιδιαίτερα, η διαφορά τους ίσως οφείλεται στο γεγονός ότι ενώ και τα δύο APIs βασίζονται στο lazy evaluation, το SQL API έχει καλύτερη διαχείριση καθώς εισάγει ολόκληρο το σχέδιο στο Catalyst Optimizer με μία εντολή. Στο DataFrame API, κάθε μετασχηματισμός εισάγεται σταδιακά, κάτι που ενδεχομένως δυσχεραίνει την εφαρμογή βελτιστοποιήσεων.

2.2 Υποερώτημα β

Στο ερώτημα αυτό, αποθηκεύσαμε το αρχικό dataset ως αρχείο τύπου parquet. Το Apache Parquet είναι ένα format αρχείου, το οποίο έχει σχεδιαστεί για να υποστηρίζει ταχύτερη επεξεργασία δεδομένων. Έπειτα, το χρησιμοποιήσαμε για να εκτελέσουμε ξανά το query 2 με το DataFrame API.

```

1 ### QUERY 2 (b)
2 crime_data.write.mode("overwrite").parquet("s3://groups-bucket-dblab-905418150721/group33/CrimeData.
  parquet")

```

Παρακάτω, φαίνεται ο τρόπος με τον οποίο διαβάσαμε τα δεδομένα από το αρχείο parquet και εκτελέσαμε ξανά το προηγούμενο query.

```

1 crime_data_parquet = spark.read.parquet("s3://groups-bucket-dblab-905418150721/group33/CrimeData.
  parquet")
2
3 def process_data_dataframe(crime_data):
4     crime_data = crime_data.withColumn('YEAR', year(to_timestamp(col('DATE OCC'), 'MM/dd/yyyy hh:mm:
5     ss a')))) \
6         .withColumn('closed', when(~col('Status').isin('IC', 'UKN'), 1).otherwise
7         (0))
8     crime_data_agg = crime_data.groupBy('YEAR', 'AREA NAME').agg(
9         count('*').alias("total_cases"),
10        sum('closed').alias('closed_cases')
11    )
12    crime_cc_rates = crime_data_agg.withColumn(
13        'closed_case_rate',
14        when(col('total_cases') > 0, col('closed_cases') / col('total_cases') * 100).otherwise(None)
15    ).drop('closed_cases').drop('total_cases')
16    window_spec = Window.partitionBy('YEAR').orderBy(desc('closed_case_rate'))
17    ranked_df = crime_cc_rates.withColumn("#", row_number().over(window_spec))
18    top3_df = ranked_df.filter(col("#") <= 3)
19    return top3_df
20
21 # DATAFRAME API BEGIN #
22 start = time.time()
23 top3_df_parquet = process_data_dataframe(crime_data_parquet)
24 top3_df_parquet_collect = top3_df_parquet.collect()
25 end = time.time()
26 print("DATAFRAME performance:", end-start)
27 print("DATAFRAME results:")

```

```
26 top3_df_parquet.show(n=top3_df_parquet.count(), truncate=False)
27 # DATAFRAME API END #
```

Η εκτέλεση του query 2 χρησιμοποιώντας τα δεδομένα από το parquet αρχείο και το DataFrame API, εκτελείται σε μόλις 4.5808703899383545 sec. Παρουσιάζει, δηλαδή, speedup **2.8778518591268054** σε σύγκριση με την εκτέλεση του query 2, χρησιμοποιώντας τα δεδομένα από τα csv αρχεία και το DataFrame API. Γεγονός που ήταν αναμενόμενο αφού η χρήση δεδομένων από parquet αρχεία είναι πολύ πιο αποδοτική λόγω της στήλη-προσανατολισμένης φύσης, της ενσωματωμένης υποστήριξης σχήματος, των βελτιστοποιήσεων φίλτρου και της αποφυγής parsing. Αντίθετα, τα csv αρχεία απαιτούν περισσότερη επεξεργασία για parsing, ανάγνωση και μετατροπή.

3 Ζητούμενο 3

Στο συγκεκριμένο ερώτημα θα χρησιμοποιήσουμε ως αναφορά τα δεδομένα απογραφής 2010 για τον πληθυσμό και εκείνα της απογραφής 2015 για το εισόδημα ανά νοικοκυριό προκειμένου να υπολογίσουμε για κάθε περιοχή του Los Angeles το μέσο ετήσιο εισόδημα ανά άτομο και την αναλογία του συνολικού αριθμού εγκλημάτων ανά άτομο. Το query εκτελέστηκε με DataFrame API.

Πρώτα απ' όλα απαιτείται να αρχικοποιήσουμε το Sedona Context, το οποίο είναι απαραίτητο για την εκτέλεση geospatial analytics χρησιμοποιώντας τη βιβλιοθήκη Apache Sedona. Τα δεδομένα εγκλημάτων φορτώνονται από δύο αρχεία csv που περιέχουν εγγραφές από το 2010 έως το παρόν και ενοποιούνται με τη μέθοδο `union`. Για να αποφευχθούν σφάλματα, αφαιρούνται εγγραφές με συντεταγμένες στο Null Island (0,0). Στη συνέχεια, προστίθεται μια νέα στήλη `geometry` που περιέχει τα γεωμετρικά σημεία των εγκλημάτων, χρησιμοποιώντας τις στήλες `LAT` και `LON`. Τα δεδομένα απογραφής πληθυσμού φορτώνονται από ένα αρχείο GeoJSON. Αφαιρούνται εγγραφές με αρνητικές τιμές πληθυσμού ή κατοικιών. Επίσης, φιλτράρονται μόνο οι εγγραφές που ανήκουν στην πόλη του Los Angeles, με βάση τη στήλη `CITY`. Τέλος τα δεδομένα εισοδήματος φορτώνονται από ένα αρχείο csv. Η στήλη που περιέχει το μέσο εισόδημα ανά νοικοκυριό καθαρίζεται αφαιρώντας τα σύμβολα \$ και ,. Η καθαρισμένη τιμή μετατρέπεται σε αριθμητικό τύπο δεδομένων (`DecimalType`).

```
1 ### QUERY 3
2 APP_NAME = "(Income/person and crimes/person) / region"
3 spark = SparkSession.builder.appName(APP_NAME).getOrCreate()
4
5 sedona = SedonaContext.create(spark)
6
7 # crime data
8 d1 = spark.read.csv(
9     "s3://initial-notebook-data-bucket-dblab-905418150721/CrimeData/
10     Crime_Data_from_2010_to_2019_20241101.csv",
11     header=True, inferSchema=True)
12 d2 = spark.read.csv(
13     "s3://initial-notebook-data-bucket-dblab-905418150721/CrimeData/
14     Crime_Data_from_2020_to_Present_20241101.csv",
15     header=True, inferSchema=True)
16 crime_data = d1.union(d2)
17
18 # remove NULL ISLAND (0,0)
19 crime_data = crime_data.filter(~((col('LAT') == 0) & (col('LON') == 0))) \
20     .select('DR_NO', 'LAT', 'LON').withColumn('geometry', ST_Point('LON', 'LAT')).drop('LON').
21     drop('LAT')
22
23 # census block data
24 geojson_path = "s3://initial-notebook-data-bucket-dblab-905418150721/2010_Census_Blocks.geojson"
25 census_blocks = sedona.read.format("geojson") \
26     .option("multiLine", "true").load(geojson_path) \
27     .selectExpr("explode(features) as features") \
28     .select("features.*")
29 census_blocks = census_blocks.select(
30     [col(f"properties.{col_name}").alias(col_name) for col_name in
31     census_blocks.schema["properties"].dataType.fieldNames()] + ["geometry"],).drop("
32     properties").drop("type")
33 census_blocks = census_blocks.filter((census_blocks['HOUSING10'] >= 0) & (census_blocks['POP_2010']
34     >= 0)).filter(F.col('CITY') == 'Los Angeles')
35
36 # median household income data
37 income_data = spark.read.csv(
38     "s3://initial-notebook-data-bucket-dblab-905418150721/LA_income_2015.csv",
39     header=True, inferSchema=True)
40
41 income_data = income_data.withColumn(
42     "median_hincome",
43     regexp_replace(col("Estimated Median Income"), "[$,]", "").cast(DecimalType())
44 ).drop('Estimated Median Income', 'Community')
```

Για την κατάλληλη επεξεργασία των δεδομένων χρησιμοποιείται η μέθοδος `ST_Within` της Sedona για να συσχετιστούν τα δεδομένα εγκλημάτων με τα δεδομένα απογραφής. Συγκεκριμένα, κάθε περιστατικό εγκλήματος αντιστοιχίζεται στην περιοχή απογραφής (`COMM`) στην οποία ανήκει. Στη συνέχεια, υπολογίζεται ο συνολικός αριθμός εγκλημάτων για κάθε περιοχή. Οι περιοχές ομαδοποιούνται με βάση τον ταχυδρομικό κώδικα (`ZCTA10`) και την περιοχή (`COMM`). Υπολογίζονται το συνολικό πλήθος κατοίκων και κατοικιών για κάθε περιοχή. Έπειτα, γίνεται σύνδεση με τα δεδομένα εισοδήματος, χρησιμοποιώντας τον ταχυδρομικό κώδικα ως κοινό πεδίο. Υπολογίζεται το συνολικό εισόδημα για κάθε περιοχή, πολλαπλασιάζοντας τον αριθμό κατοικιών με το μέσο εισόδημα ανά νοικοκυριό. Οι περιοχές ομαδοποιούνται με βάση τη στήλη `COMM`, και υπολογίζεται το συνολικό εισόδημα και ο πληθυσμός

για κάθε περιοχή. Στη συνέχεια, συνδέονται τα δεδομένα εγκλημάτων με αυτά τα σύνολα, δημιουργώντας έναν ενιαίο πίνακα που περιλαμβάνει τα δεδομένα εγκλημάτων, πληθυσμού και εισοδήματος για κάθε περιοχή (έγινε χρήση RIGHT JOIN, έτσι ώστε να συμπεριλάβουμε περιοχές χωρίς κανένα έγκλημα). Τέλος για κάθε περιοχή, υπολογίζεται το μέσο εισόδημα ανά άτομο και η αναλογία εγκλημάτων ανά άτομο, σε περιπτώσεις όπου ο πληθυσμός είναι μηδενικός, ορίζεται η τιμή σε μηδέν.

```

1 # DATAFRAME API BEGIN #
2 start = time.time()
3 # crime data
4 comm_join_crimes = crime_data.join(census_blocks, ST_Within(crime_data['geometry'], census_blocks['
    geometry']))
5 count_comm_crimes = comm_join_crimes.groupby('COMM').agg(F.count('*').alias('crimes_count'))
6 #income data
7 zip_comm_houses_pop = census_blocks.groupBy('ZCTA10', 'COMM').agg(
8     sum(col('POP_2010')).alias('population'),
9     sum(col('HOUSING10')).alias('houses'),
10 ).select('ZCTA10', 'COMM', 'population', 'houses')
11 # joined with median household income
12 zip_comm_houses_pop_hincome = zip_comm_houses_pop.join(
13     income_data,
14     income_data['Zip Code'] == zip_comm_houses_pop['ZCTA10']
15 ).drop('Zip Code')
16 # zip_total_income = houses * median_hincome per (zip, comm)
17 zip_comm_total_income = zip_comm_houses_pop_hincome.withColumn(
18     "zip_total_income",
19     col('median_hincome') * col('houses')
20 )
21 # total_population per COMM
22 comm_total_population_total_income = zip_comm_total_income.groupBy('COMM').agg(
23     sum('population').alias('total_population'),
24     sum('zip_total_income').alias('total_income')
25 )
26 # join crime and income data
27 comm_crime_income = comm_total_population_total_income.join(
28     count_comm_crimes,
29     on='COMM',
30     how='right'
31 )
32 comm_mincome_person = comm_crime_income.withColumn(
33     'average_income_per_person',
34     when(col('total_population') > 0, col('total_income') / col('total_population'))
35     .otherwise(0)).withColumn(
36     'crime_rate_per_person',
37     when(col('total_population') > 0, col('crimes_count') / col('total_population'))
38     .otherwise(0)
39 ).drop('total_income').drop('crimes_count').drop('total_population')
40 comm_mincome_person_collect = comm_mincome_person.collect()
41 end = time.time()
42 print("DATAFRAME performance:", end-start)
43 comm_mincome_person.show(n=comm_mincome_person.count(), truncate=False)
44 # DATAFRAME API END #

```

Παρακάτω φαίνεται η έξοδος του ερωτήματος (οι πρώτες 20 σειρές) το οποίο εκτελείται σε 29.66372036933899 sec.

COMM	average_income_per_person	crime_rate_per_person
San Pedro	24244.561240	0.7751021460858706
South Park	6943.255564	0.8862840357021805
Vernon Central	6624.199012	0.7288347583605012
Florence-Firestone	8079.274829	1.1133874146386178
Studio City	44049.852630	0.9295754238516157
Valley Village	28191.368349	0.6809068681929318
Sherman Oaks	37767.445674	0.7903687241383545
Valley Glen	18271.771022	0.5546786523216308
Park La Brea	36619.900017	1.025207944321847
Miracle Mile	38834.846111	0.8820452139253908
Hollywood	25648.050759	1.5511440107671601
Hancock Park	21538.793405	0.8983094426946069
Melrose	21688.494261	0.909787014117455
North Hills	16224.486386	0.655059866962306
Panorama City	10221.126106	0.6498237536867851
Mission Hills	17158.531061	0.6565328178543186
Arleta	12110.779170	0.4423287504562599

Για να δοκιμάσουμε όλους τους δυνατούς συνδυασμούς στις στρατηγικές join μεταξύ των BROADCAST, MERGE, SHUFFLE_HASH, SHUFFLE_REPLICATE_NL για τα τρία join του query αλλάζουμε τον κώδικα ως εξής:

```

1  ### QUERY 3 (join strategies)
2  def execute_query_with_join_strategy(strategy1, strategy2, strategy3):
3      start = time.time()
4
5      # Ρυθμίσεις για το Spark SQL για κάθε join
6      spark.conf.set("spark.sql.autoBroadcastJoinThreshold", -1) # Απενεργοποίηση default broadcast
7      spark.conf.set("spark.sql.join.preferSortMergeJoin", "true" if strategy1 == "MERGE" else "false"
8      )
9      spark.conf.set("spark.sql.join.preferBroadcastJoin", "true" if strategy1 == "BROADCAST" else "
10      false")
11
12      # crime_data join census_blocks (strategy1)
13      comm_join_crimes = crime_data.join(
14          census_blocks.hint(strategy1),
15          ST_Within(crime_data['geometry'], census_blocks['geometry'])
16      )
17      count_comm_crimes = comm_join_crimes.groupby('COMM').agg(F.count('*').alias('crimes_count'))
18
19      # zip_comm_houses_pop_hincome join income_data (strategy2)
20      zip_comm_houses_pop = census_blocks.groupBy('ZCTA10', 'COMM').agg(
21          sum(col('POP_2010')).alias('population'),
22          sum(col('HOUSING10')).alias('houses'),
23      ).select('ZCTA10', 'COMM', 'population', 'houses')
24
25      zip_comm_houses_pop_hincome = zip_comm_houses_pop.join(
26          income_data.hint(strategy2),
27          income_data['Zip Code'] == zip_comm_houses_pop['ZCTA10']
28      ).drop('Zip Code')
29
30      # zip_comm_total_income (strategy3)
31      zip_comm_total_income = zip_comm_houses_pop_hincome.withColumn(
32          "zip_total_income",
33          col('median_hincome') * col('houses')
34      )
35
36      comm_total_population_total_income = zip_comm_total_income.groupBy('COMM').agg(
37          sum('population').alias('total_population'),
38          sum('zip_total_income').alias('total_income')
39      )
40
41      comm_crime_income = comm_total_population_total_income.join(
42          count_comm_crimes.hint(strategy3),
43          on='COMM',
44          how='outer'
45      )
46
47      comm_mincome_person = comm_crime_income.withColumn(
48          'average_income_per_person',
49          when(col('total_population') > 0, col('total_income') / col('total_population'))
50          .otherwise(0)
51      ).withColumn(
52          'crime_rate_per_person',
53          when(col('total_population') > 0, col('crimes_count') / col('total_population'))
54          .otherwise(0)
55      ).drop('total_income').drop('crimes_count').drop('total_population')
56      comm_mincome_person_collect = comm_mincome_person.collect()
57
58      end = time.time()
59      execution_time = end - start
60
61      return execution_time
62
63  join_strategies = ["BROADCAST", "MERGE", "SHUFFLE_HASH", "SHUFFLE_REPLICATE_NL"]
64  strategy_combinations = list(itertools.product(join_strategies, repeat=3))
65
66  results = []
67
68  for strat1, strat2, strat3 in strategy_combinations:
69      exec_time = execute_query_with_join_strategy(strat1, strat2, strat3)
70      results.append((strat1, strat2, strat3, exec_time))
71
72  results_df = spark.createDataFrame(results, ["Strategy1", "Strategy2", "Strategy3", "ExecutionTime"
73  ])
74  results_df_sorted = results_df.orderBy("ExecutionTime", ascending=True)

```

```
72 results_df_sorted.show(n=results_df_sorted.count(), truncate=False)
```

Παρακάτω φαίνεται ο χρόνος εκτέλεσης για κάθε συνδυασμό join strategies.

Strategy1	Strategy2	Strategy3	ExecutionTime
SHUFFLE_REPLICATE_NL	SHUFFLE_HASH	BROADCAST	13.472944498062134
BROADCAST	SHUFFLE_REPLICATE_NL	SHUFFLE_REPLICATE_NL	13.713311672210693
SHUFFLE_REPLICATE_NL	SHUFFLE_REPLICATE_NL	MERGE	13.72116470336914
SHUFFLE_HASH	SHUFFLE_REPLICATE_NL	SHUFFLE_REPLICATE_NL	13.74719762802124
SHUFFLE_HASH	SHUFFLE_HASH	BROADCAST	13.748756647109985
SHUFFLE_REPLICATE_NL	MERGE	BROADCAST	14.009728908538818
SHUFFLE_REPLICATE_NL	BROADCAST	BROADCAST	14.104211568832397
SHUFFLE_HASH	SHUFFLE_HASH	SHUFFLE_HASH	14.124082088470459
SHUFFLE_HASH	SHUFFLE_HASH	MERGE	14.143094301223755
SHUFFLE_HASH	MERGE	MERGE	14.165607452392578
BROADCAST	SHUFFLE_HASH	MERGE	14.181203842163086
SHUFFLE_REPLICATE_NL	BROADCAST	SHUFFLE_HASH	14.28413438796997
SHUFFLE_REPLICATE_NL	MERGE	SHUFFLE_HASH	14.296715497970581
SHUFFLE_HASH	SHUFFLE_REPLICATE_NL	SHUFFLE_HASH	14.303778886795044
MERGE	MERGE	BROADCAST	14.306119918823242
BROADCAST	SHUFFLE_REPLICATE_NL	MERGE	14.30896806716919
SHUFFLE_REPLICATE_NL	SHUFFLE_REPLICATE_NL	SHUFFLE_REPLICATE_NL	14.312310218811035
SHUFFLE_HASH	SHUFFLE_REPLICATE_NL	BROADCAST	14.338404417037964
SHUFFLE_HASH	BROADCAST	SHUFFLE_REPLICATE_NL	14.344301700592041
SHUFFLE_REPLICATE_NL	MERGE	MERGE	14.345511674880981
MERGE	SHUFFLE_HASH	BROADCAST	14.434895992279053
SHUFFLE_REPLICATE_NL	MERGE	SHUFFLE_REPLICATE_NL	14.442206621170044
SHUFFLE_REPLICATE_NL	BROADCAST	SHUFFLE_REPLICATE_NL	14.490171670913696
MERGE	SHUFFLE_HASH	SHUFFLE_HASH	14.524138689041138
MERGE	SHUFFLE_REPLICATE_NL	MERGE	14.545830965042114
SHUFFLE_HASH	SHUFFLE_REPLICATE_NL	MERGE	14.566548347473145
MERGE	SHUFFLE_REPLICATE_NL	SHUFFLE_HASH	14.573344707489014
MERGE	SHUFFLE_HASH	MERGE	14.653209447860718
BROADCAST	SHUFFLE_HASH	SHUFFLE_REPLICATE_NL	14.756784200668335
SHUFFLE_HASH	SHUFFLE_HASH	SHUFFLE_REPLICATE_NL	14.820102214813232
MERGE	SHUFFLE_HASH	SHUFFLE_REPLICATE_NL	14.938286542892456
BROADCAST	SHUFFLE_HASH	BROADCAST	14.974244832992554
SHUFFLE_HASH	BROADCAST	BROADCAST	14.982024669647217
SHUFFLE_HASH	MERGE	BROADCAST	15.060087203979492
SHUFFLE_REPLICATE_NL	SHUFFLE_HASH	SHUFFLE_HASH	15.076187133789062
SHUFFLE_HASH	MERGE	SHUFFLE_HASH	15.079508543014526
SHUFFLE_REPLICATE_NL	SHUFFLE_REPLICATE_NL	BROADCAST	15.351330280303955
SHUFFLE_REPLICATE_NL	SHUFFLE_HASH	SHUFFLE_REPLICATE_NL	15.469627141952515
SHUFFLE_HASH	BROADCAST	SHUFFLE_HASH	15.482222080230713
BROADCAST	BROADCAST	SHUFFLE_HASH	15.497687578201294
MERGE	MERGE	MERGE	15.682572364807129
BROADCAST	BROADCAST	MERGE	15.800673723220825
SHUFFLE_HASH	MERGE	SHUFFLE_REPLICATE_NL	15.830183982849121
MERGE	MERGE	SHUFFLE_HASH	15.844552516937256
MERGE	BROADCAST	BROADCAST	15.85817003250122
MERGE	BROADCAST	SHUFFLE_HASH	15.874963998794556
MERGE	SHUFFLE_REPLICATE_NL	SHUFFLE_REPLICATE_NL	15.993769884109497
SHUFFLE_REPLICATE_NL	SHUFFLE_REPLICATE_NL	SHUFFLE_HASH	16.042428493499756
BROADCAST	BROADCAST	SHUFFLE_REPLICATE_NL	16.231744050979614
BROADCAST	MERGE	MERGE	16.40001654624939
BROADCAST	MERGE	SHUFFLE_REPLICATE_NL	16.669507026672363
SHUFFLE_REPLICATE_NL	BROADCAST	MERGE	17.152887105941772
SHUFFLE_HASH	BROADCAST	MERGE	17.63453984260559
BROADCAST	BROADCAST	BROADCAST	17.746803045272827
MERGE	BROADCAST	MERGE	17.746988534927368
BROADCAST	MERGE	BROADCAST	18.210407495498657
BROADCAST	SHUFFLE_HASH	SHUFFLE_HASH	19.552934885025024
MERGE	BROADCAST	SHUFFLE_REPLICATE_NL	19.877957344055176
MERGE	SHUFFLE_REPLICATE_NL	BROADCAST	20.444053649902344
MERGE	MERGE	SHUFFLE_REPLICATE_NL	20.667855501174927
BROADCAST	SHUFFLE_REPLICATE_NL	BROADCAST	20.7849862575531
BROADCAST	MERGE	SHUFFLE_HASH	21.747425079345703
BROADCAST	SHUFFLE_REPLICATE_NL	SHUFFLE_HASH	23.430171728134155
SHUFFLE_REPLICATE_NL	SHUFFLE_HASH	MERGE	24.037864923477173

Γενικότερα οι οι στρατηγικές join καθορίζουν πώς δύο σύνολα δεδομένων συνδυάζονται. Οι κύριες στρατηγικές είναι:

- **Broadcast Join:** Ένα μικρότερο dataset αντιγράφεται και διαμοιράζεται σε όλους τους κόμβους. Είναι ιδανικό για μικρά datasets, καθώς μειώνει την ανάγκη για data shuffling.
- **Merge Join (Sort-Merge Join):** Τα datasets ταξινομούνται πρώτα και στη συνέχεια συνδυάζονται. Αυτή η στρατηγική είναι αποτελεσματική για μεγάλα datasets που είναι ήδη ταξινομημένα.
- **Shuffle Hash Join:** Τα δεδομένα χωρίζονται σε partitions βάσει του hash key και στη συνέχεια γίνεται το join. Είναι κατάλληλο για μεγάλου μεγέθους datasets αλλά απαιτεί shuffling.
- **Shuffle Replicate NL (Nested Loop Join):** Το Spark χρησιμοποιεί shuffling στους πίνακες, ενώ μετά αντιγράφεται ένας από τους δύο προκύπτοντες πίνακες και το αποτέλεσμα είναι ένα καρτεσιανό γινόμενο των δύο datasets.

Από τα αποτελέσματα που παρουσιάζονται στον παραπάνω πίνακα έχουμε τις εξής παρατηρήσεις:

- Ο συνδυασμός SHUFFLE_REPLICATE_NL, SHUFFLE_HASH, και BROADCAST είναι ο γρηγορότερος με χρόνο εκτέλεσης 13.47 δευτερόλεπτα. Αυτό οφείλεται στη χρήση του BROADCAST που ελαχιστοποιεί το shuffling για το τρίτο join.
- Ο συνδυασμός BROADCAST, SHUFFLE_HASH, και SHUFFLE_REPLICATE_NL έχει τον μεγαλύτερο χρόνο εκτέλεσης (23.43 δευτερόλεπτα), λόγω της χρήσης του S SHUFFLE_REPLICATE_NL, που είναι αργό για μεγάλα δεδομένα.
- Η χρήση του BROADCAST γενικά μειώνει το χρόνο εκτέλεσης όταν χρησιμοποιείται σε κρίσιμα σημεία.
- Οι στρατηγικές MERGE και SHUFFLE_HASH δείχνουν καλή ισορροπία απόδοσης.
- Ο SHUFFLE_REPLICATE_NL αυξάνει σημαντικά το χρόνο εκτέλεσης όταν χρησιμοποιείται επαναλαμβανόμενα.
- Συνδυασμοί που περιλαμβάνουν το BROADCAST στην τελευταία ή πρώτη φάση τείνουν να είναι γρηγορότεροι.
- Όταν το MERGE χρησιμοποιείται σε τουλάχιστον μία φάση, μειώνεται το κόστος ταξινόμησης και shuffling.
- Η επιλογή των στρατηγικών join πρέπει να γίνεται βάσει του μεγέθους των δεδομένων και της σχέσης μεταξύ τους.
- Το BROADCAST είναι ιδανικό για μικρά datasets, ενώ το MERGE ή το SHUFFLE_HASH για μεγάλα.

4 Ζητούμενο 4

Σε αυτό το ερώτημα καλούμαστε να βρούμε το φυλετικό προφίλ των καταγεγραμμένων θυμάτων εγκλημάτων στο Los Angeles για το έτος 2015 τόσο στις τρεις περιοχές με το υψηλότερο κατά κεφαλήν εισόδημα όσο και στις τρεις περιοχές με το χαμηλότερο κατά κεφαλήν εισόδημα. Τα αποτελέσματα ζητούνται σε φθίνουσα σειρά πλήθους θυμάτων ανά φυλετική κατηγορία. Το query εκτελέστηκε με DataFrame API.

Αρχικά όπως και στο προηγούμενο ερώτημα γίνεται φόρτωση των δεδομένων εγκλημάτων από ένα αρχείο csv (χρησιμοποιούμε μόνο αυτό που περιέχει δεδομένα για τα έτη 2010-2019 αφού τα αποτελέσματα μας αφορούν το έτος 2015). Τα δεδομένα φιλτράρονται για να απομακρυνθούν οι εγγραφές με συντεταγμένες (0,0). Επιπλέον, για κάθε εγγραφή δημιουργείται μια γεωμετρική αναπαράσταση (geometry) των συντεταγμένων με τη βοήθεια της βιβλιοθήκης Apache Sedona. Έπειτα τα δεδομένα φιλτράρονται ώστε να διατηρηθούν μόνο τα εγκλήματα που συνέβησαν το 2015, και έχουν συμπληρωμένη τη στήλη Vict Descent που περιγράφει τη φυλετική καταγωγή του θύματος. Στην συνέχεια ακολουθούμε ακριβώς την ίδια διαδικασία με το ερώτημα τρία για να εξάγουμε το μέσο εισόδημα ανά άτομο σε κάθε περιοχή.

```
1 ### Query 4
2 APP_NAME = "Crime Victime Race Analysis"
3 SPARK_EXECUTORS = 2
4 spark = SparkSession.builder.appName(APP_NAME).config("spark.executor.instances", SPARK_EXECUTORS).
5   getOrCreate()
6 sedona = SedonaContext.create(spark)
7 # crime data
8 crime_data = spark.read.csv(
9   "s3://initial-notebook-data-bucket-dblab-905418150721/CrimeData/
10   Crime_Data_from_2010_to_2019_20241101.csv",
11   header=True, inferSchema=True
12 )
13 # remove NULL ISLAND (0,0)
14 crime_data = crime_data.filter(~((F.col('LAT') == 0) & (F.col('LON') == 0))) \
15   .select('DR_NO', 'LAT', 'LON', 'Vict Descent', 'Date Rptd') \
16   .withColumn('geometry', ST_Point('LON', 'LAT')) \
17   .drop('LON').drop('LAT')
18
19 # filter crimes for the year 2015
20 crime_data = crime_data.withColumn("Date Rptd", to_timestamp(F.col("Date Rptd"), "MM/dd/yyyy hh:mm:
21   ss a"))
22 crimes_2015 = crime_data.filter(
23   (F.year(F.col("Date Rptd")) == 2015) & (F.col("Vict Descent").isNotNull())
24 )
25 # median household income data
26 income_data = spark.read.csv(
27   "s3://initial-notebook-data-bucket-dblab-905418150721//LA_income_2015.csv",
28   header=True, inferSchema=True
29 )
30 income_data = income_data.withColumn(
31   "median_income",
32   F.regexp_replace(F.col("Estimated Median Income"), "[,]", "").cast(DecimalType())
33 ).drop('Estimated Median Income', 'Community')
34
35 # census block data
36 geojson_path = "s3://initial-notebook-data-bucket-dblab-905418150721/2010_Census_Blocks.geojson"
37 census_blocks = sedona.read.format("geojson") \
38   .option("multiLine", "true").load(geojson_path) \
39   .selectExpr("explode(features) as features") \
40   .select("features.*")
41 census_blocks = census_blocks.select(
42   [col(f"properties.{col_name}").alias(col_name) for col_name in
43     census_blocks.schema["properties"].dataType.fieldNames()] + ["geometry"],).drop("
44   properties").drop("type")
45 census_blocks = census_blocks.filter(F.col('CITY') == 'Los Angeles')
46
47 # race and ethnicity legend
48 descent_legend = spark.read.csv(
49   "s3://initial-notebook-data-bucket-dblab-905418150721/RE_codes.csv",
50   header=True, inferSchema=True
51 )
52
53 ### Query 3
54 comm_join_crimes = crime_data.join(census_blocks, ST_Within(crime_data['geometry'], census_blocks['
55   geometry']))
```

```

54 count_comm_crimes = comm_join_crimes.groupby('COMM').agg(F.count('*').alias('crimes_count'))
55
56 zip_comm_houses_pop = census_blocks.groupBy('ZCTA10', 'COMM').agg(
57     sum(col('POP_2010')).alias('population'),
58     sum(col('HOUSING10')).alias('houses'),
59 ).select('ZCTA10', 'COMM', 'population', 'houses')
60
61 zip_comm_houses_pop_hincome = zip_comm_houses_pop.join(
62     income_data,
63     income_data['Zip Code'] == zip_comm_houses_pop['ZCTA10']
64 ).drop('Zip Code')
65
66 zip_comm_total_income = zip_comm_houses_pop_hincome.withColumn(
67     "zip_total_income",
68     col('median_income') * col('houses')
69 )
70
71 comm_total_population_total_income = zip_comm_total_income.groupBy('COMM').agg(
72     sum('population').alias('total_population'),
73     sum('zip_total_income').alias('total_income')
74 )
75
76 comm_crime_income = comm_total_population_total_income.join(
77     count_comm_crimes,
78     on='COMM',
79     how='right'
80 )
81 comm_mincome_person = comm_crime_income.withColumn(
82     'average_income_per_person',
83     when(col('total_population') > 0, col('total_income') / col('total_population'))
84     .otherwise(0)).withColumn(
85     'crime_rate_per_person',
86     when(col('total_population') > 0, col('crimes_count') / col('total_population'))
87     .otherwise(0)
88 ).drop('total_income').drop('crimes_count').drop('total_population')

```

Με τη χρήση ενός παραθύρου (window function), υπολογίζεται ένας μοναδικός αριθμός σειράς για κάθε κοινότητα με βάση το μέσο εισόδημα ανά άτομο. Έτσι, εντοπίζονται οι τρεις περιοχές με το υψηλότερο και χαμηλότερο κατά κεφαλήν εισόδημα. Τα δεδομένα εγκλημάτων συνδέονται με τα census blocks, χρησιμοποιώντας τη γεωχωρική συνάρτηση ST_Within της Apache Sedona, ώστε να αποδοθεί κάθε έγκλημα στην αντίστοιχη περιοχή. Στη συνέχεια, πραγματοποιείται αντιστοίχιση με τον πίνακα descent_legend, ώστε να μετατραπούν οι κωδικοί φυλετικής καταγωγής σε πλήρη περιγραφή. Τα εγκλήματα για τις τρεις περιοχές με το υψηλότερο και χαμηλότερο εισόδημα ομαδοποιούνται ανά φυλετική καταγωγή (descent_description), υπολογίζεται ο αριθμός των θυμάτων για κάθε φυλετική ομάδα και τα αποτελέσματα εμφανίζονται ταξινομημένα σε φθίνουσα σειρά σε δύο πίνακες.

```

1 def query4(crimes_2015, income_data, census_blocks, descent_legend, comm_mincome_person):
2     window_spec = Window.orderBy("average_income_per_person")
3     comm_mincome_person = comm_mincome_person.withColumn("row_num", row_number().over(window_spec))
4     # with row number
5     comm_bottom_3 = comm_mincome_person.filter(comm_mincome_person.row_num <= 3).select('COMM')
6     comm_top_3 = comm_mincome_person.filter(comm_mincome_person.row_num > comm_mincome_person.count() - 3).select("COMM")
7
8     # match crimes to zip codes
9     crimes_with_zipcodes = crimes_2015.join(
10         census_blocks,
11         ST_Within(crime_data['geometry'], census_blocks['geometry'])
12     )
13
14     # join with descent legend for descriptions
15     crimes_with_descriptions = crimes_with_zipcodes.join(
16         descent_legend,
17         crimes_with_zipcodes["Vict Descent"] == descent_legend["Vict Descent"],
18         how="left"
19     ).select(
20         crimes_with_zipcodes["*"],
21         descent_legend["Vict Descent Full"].alias("descent_description")
22     )
23
24     top_crimes = crimes_with_descriptions.join(comm_top_3, "COMM", "inner")
25     bottom_crimes = crimes_with_descriptions.join(comm_bottom_3, "COMM", "inner")
26
27     # group and count victims by descent description for each group
28     top_victim_count = top_crimes.groupBy("descent_description") \
29         .agg(F.count("*").alias("victim_count")) \

```

```

29     .orderBy(F.desc("victim_count"))
30
31     bottom_victim_count = bottom_crimes.groupBy("descent_description") \
32     .agg(F.count("*").alias("victim_count")) \
33     .orderBy(F.desc("victim_count"))
34
35     # display results
36     print("Top COMM Victim Count:")
37     top_victim_count.show()
38
39     print("Bottom COMM Victim Count:")
40     bottom_victim_count.show()

```

Σημείωση: Θα μπορούσαμε να εκτελέσουμε την πράξη UNION μεταξύ των 3 περιοχών με χαμηλότερο εισόδημα και των 3 με μεγαλύτερο, ώστε να γίνει ένα JOIN με το dataset εγκλημάτων αντί για δύο. Θέλαμε, ωστόσο, να τα έχουμε σε διαφορετικά DataFrames.

Τα αποτελέσματα του ερωτήματος φαίνονται παρακάτω.

```

1 Top COMM Victim Count:
2 +-----+-----+
3 | descent_description|victim_count|
4 +-----+-----+
5 |                White|          682|
6 |                Other|           81|
7 |Hispanic/Latin/Me...|           75|
8 |                Unknown|          50|
9 |                Black|           46|
10 |            Other Asian|           22|
11 |                Chinese|            1|
12 |American Indian/A...|            1|
13 +-----+-----+
14
15 Bottom COMM Victim Count:
16 +-----+-----+
17 | descent_description|victim_count|
18 +-----+-----+
19 |Hispanic/Latin/Me...|         3314|
20 |                Black|        1152|
21 |                White|         435|
22 |                Other|         250|
23 |            Other Asian|         136|
24 |                Unknown|          31|
25 |American Indian/A...|          22|
26 |                Chinese|           4|
27 |                Korean|           4|
28 |                Filipino|           3|
29 |            AsianIndian|            1|
30 |                Guamanian|            1|
31 +-----+-----+

```

Τέλος μας ζητείται να εκτελέσουμε το ερώτημα με διαφορετικά configurations (διατηρώντας σταθερό και ίσο με 2 το πλήθος των executors) έτσι ώστε να δούμε πως αυτή ανταποκρίνεται στην κλιμάκωση του συνόλου των υπολογιστικών πόρων χρήσης. Τα αποτελέσματα φαίνονται στον παρακάτω πίνακα.

Executor Cores	Executor Memory	Processing Time (s)
1	2GB	121.39
2	4GB	106.00
4	8GB	100.25

Τα αποτελέσματα δείχνουν ότι η αύξηση των διαθέσιμων υπολογιστικών πόρων (πυρήνες και μνήμη ανά executor) βελτιώνει την απόδοση της επεξεργασίας του ερωτήματος. Συγκεκριμένα, παρατηρείται ότι καθώς αυξάνεται ο αριθμός των πυρήνων ανά executor και η διαθέσιμη μνήμη, ο χρόνος εκτέλεσης μειώνεται. Η διαφορά μεταξύ του πρώτου και του τελευταίου configuration (με 1 πυρήνα και 2 GB μνήμης έναντι 4 πυρήνων και 8 GB μνήμης) είναι σημαντική, μειώνοντας τον χρόνο εκτέλεσης από 121.39 σε 100.25 δευτερόλεπτα. Αυτό υποδεικνύει ότι η εφαρμογή είναι σε θέση να αξιοποιήσει την αύξηση πόρων, αν και η βελτίωση από το δεύτερο στο τρίτο configuration είναι μικρότερη, υποδηλώνοντας πιθανό φαινόμενο κορεσμού ή μειωμένη απόδοση κλιμάκωσης.

5 Ζητούμενο 5

Σε αυτό το ερώτημα, καλούμαστε να υπολογίσουμε ανά αστυνομικό τμήμα, τον αριθμό εγκλημάτων που έλαβαν χώρα πλησιέστερα σε αυτό, καθώς και την μέση απόσταση του από τις τοποθεσίες όπου σημειώθηκαν τα συγκεκριμένα περιστατικά.

Αρχικά, αντλούμε τα δεδομένα για τα εγκλήματα από το crime dataset, όπως και προηγουμένως φροντίζουμε να αφαιρέσουμε εγγραφές που αναφέρονται στο Null Island (0,0). Τα δεδομένα για τις τοποθεσίες των αστυνομικών τμημάτων, αντλούνται από το αντίστοιχο dataset.

```
1 # crime data
2 d1 = spark.read.csv(
3     "s3://initial-notebook-data-bucket-dblab-905418150721/CrimeData/
4     Crime_Data_from_2010_to_2019_20241101.csv",
5     header=True, inferSchema=True)
6 d2 = spark.read.csv(
7     "s3://initial-notebook-data-bucket-dblab-905418150721/CrimeData/
8     Crime_Data_from_2020_to_Present_20241101.csv",
9     header=True, inferSchema=True)
10 crime_data = d1.union(d2)
11
12 # remove NULL ISLAND (0,0)
13 crime_data = crime_data.filter(~((col('LAT') == 0) & (col('LON') == 0)))
14 crime_data = crime_data.withColumn("crime_geometry", ST_Point(F.col("LON"), F.col("LAT")))
15 crime_data = crime_data.select(col('crime_geometry'), col('AREA '), col('DR_NO'))
16
17 # police stations data
18 police_stations = spark.read.csv(
19     "s3://initial-notebook-data-bucket-dblab-905418150721/LA_Police_Stations.csv",
20     header=True, inferSchema=True).withColumnRenamed("X", "PS_LON").withColumnRenamed("Y", "PS_LAT")
21 police_stations = police_stations.withColumn("ps_geometry", ST_Point(F.col("PS_LON"), F.col("PS_LAT")
22     )),.select(col('ps_geometry'), col('DIVISION'))
```

Για την παραγωγή των αποτελεσμάτων, αρχικά υπολογίζονται με crossjoin οι αποστάσεις όλων των εγκλημάτων από όλα τα αστυνομικά τμήματα. Για τον υπολογισμό των αποστάσεων, έγινε χρήση της συνάρτησης ST_DistanceSphere, η οποία υπολογίζει την απόσταση haversine μεταξύ δύο σημείων. Ωστόσο, επειδή τα σημεία είναι σχετικά κοντά μεταξύ τους, θα μπορούσαμε να χρησιμοποιήσουμε και την συνάρτηση ST_Distance, η οποία υπολογίζει την ευκλείδεια απόσταση μεταξύ 2 σημείων. Στο βήμα αυτό, θέλουμε για κάθε έγκλημα, να βρούμε το κοντινότερο αστυνομικό τμήμα. Ορίζουμε ένα Window Function, το οποίο κάνει partition τα δεδομένα που παρήχθησαν στο προηγούμενο βήμα με βάση την στήλη DR_NO. Το DR_NO οποίο είναι μοναδικό για κάθε σημείο στο dataset μας. Για κάθε ένα από αυτά τα partitions, κρατάμε το row με την ελάχιστη απόσταση μεταξύ αστυνομικού τμήματος και σημείου εγκλήματος. Τελικά, ομαδοποιούμε τα δεδομένα, που παρήχθησαν κρατώντας τις ελάχιστες αποστάσεις, με βάση το DIVISION (δηλαδή το αστυνομικό τμήμα) και βρίσκουμε τον μέσο όρο απόστασης ανά αστυνομικό τμήμα.

```
1 def query5(crime_data, police_stations):
2     joined_df = crime_data.crossJoin(police_stations).withColumn(
3         "distance_km",
4         ST_DistanceSphere(F.col("crime_geometry"), F.col("ps_geometry")) / 1000
5     )
6
7     window_spec = Window.partitionBy("DR_NO").orderBy(asc("distance_km"))
8     joined_df_with_rank = joined_df.withColumn("rank", F.row_number().over(window_spec))
9     nearest_station_df = joined_df_with_rank.filter(F.col("rank") == 1).select(
10         "DR_NO", "DIVISION", "distance_km"
11     )
12
13     result_df = nearest_station_df.groupBy("DIVISION").agg(
14         F.count("*").alias("#"),
15         F.mean("distance_km").alias("average_distance")
16     ).orderBy(F.col("#").desc())
17
18     result_df.show(n=result_df.count(), truncate=False)
```

Παρακάτω φαίνονται τα αποτελέσματα του query 5.

```
1 +-----+-----+-----+
2 |DIVISION      |#      |average_distance |
3 +-----+-----+-----+
4 |HOLLYWOOD     |224340 |2.076263960178719 |
5 |VAN NUYS      |210134 |2.9533697428197816 |
6 |SOUTHWEST     |188901 |2.191398805780884 |
7 |WILSHIRE      |185996 |2.5926655329787645 |
```



```

8 |77TH STREET      |171827|1.716544971970096 |
9 |OLYMPIC          |170897|1.7236036971780955|
10|NORTH HOLLYWOOD  |167854|2.643006094141565 |
11|PACIFIC          |161359|3.8500706553079165|
12|CENTRAL          |153871|0.9924764374568908|
13|RAMPART          |152736|1.5345341879190117|
14|SOUTHEAST        |152176|2.4218662158881807|
15|WEST VALLEY      |138643|3.035671216314069 |
16|TOPANGA          |138217|3.2969548417555528|
17|FOOTHILL         |134896|4.250921708424983 |
18|HARBOR           |126747|3.7025615993565206|
19|HOLLENBECK       |115837|2.680181237706819 |
20|WEST LOS ANGELES|115781|2.7924572890341217|
21|NEWTON           |111110|1.6346357397097429|
22|NORTHEAST        |108109|3.6236655246040868|
23|MISSION          |103355|3.6909426142785966|
24|DEVONSHIRE       |77094 |2.824765412800822 |
25+-----+-----+-----+

```

Τέλος μας ζητείται να εκτελέσουμε το ερώτημα με διαφορετικά configurations χρησιμοποιώντας σταθερούς συνολικούς πόρους 8 cores και 16GB μνήμης. Τα αποτελέσματα φαίνονται στον παρακάτω πίνακα.

Executors	Executor Cores	Executor Memory	Processing Time (s)
2	4	8GB	19.10
4	2	4GB	37.56
8	1	2GB	41.06

Τα αποτελέσματα δείχνουν, πως ο χρόνος εκτέλεσης μειώνεται δραστικά όταν έχουμε περισσότερη μνήμη σε κάθε executor. Αυτό οφείλεται πιθανότατα στο γεγονός πως υπάρχει αρκετή μνήμη ώστε οι executors, να αποθηκεύουν ενδιάμεσα αποτελέσματα στην κύρια μνήμη. Σε περίπτωση που δεν επαρκεί η κύρια μνήμη, για τα ενδιάμεσα αποτελέσματα θα πρέπει να γίνουν αποθηκεύσεις στον δίσκο. Το πρώτο configuration παρουσιάζει την καλύτερη επίδοση, γιατί έχει αρκετή μνήμη για να αποφύγουμε I/O operations με τον δίσκο, αλλά και επειδή έχει αρκετό αριθμό cores ώστε να υπάρχει ένας ικανοποιητικός βαθμός παραλληλίας.