

ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ ΣΧΟΛΗ ΗΛΕΚΤΡΟΛΟΓΩΝ ΜΗΧΑΝΙΚΩΝ ΚΑΙ ΜΗΧΑΝΙΚΩΝ ΥΠΟΛΟΓΙΣΤΩΝ

ΤΟΜΕΑΣ ΤΕΧΝΟΛΟΓΙΑΣ ΠΛΗΡΟΦΟΡΙΚΗΣ ΚΑΙ ΥΠΟΛΟΓΙΣΤΩΝ ΕΡΓΑΣΤΗΡΙΟ ΒΑΣΕΩΝ ΓΝΩΣΕΩΝ ΚΑΙ ΔΕΔΟΜΕΝΩΝ

Ακαδημαϊκό έτος 2023-24, 9ο Εξάμηνο

Εξαμηνιαία εργασία στα Προχωρημένα συστήματα Βάσεων Δεδομένων

Ομάδα 24: Ευάγγελος Παπαμήτσος, Χρήστος Μεντζίνης

Η εργασία πραγματοποιήθηκε τοπικά σε έναν κόμβο λογω μη λήψης πόρων στον ωκεανό και αποτυχίας δημιοργίας του κατάλληλου distributed περιβάλλοντος με VMs, docker...

Δημιουργούμε ένα DataFrame από το σύνολο δεδομένων και αφού κάνουμε τις ζητούμενες μετατροπές πάιρνουμε τις παρακάτω πληροφορίες:

Total Rows: 2979324

DR_NO: integer (nullable = true)	Date Rptd: date (nullable = true)
DATE OCC: date (nullable = true)	TIME OCC: integer (nullable = true)
AREA: integer (nullable = true)	AREA NAME: string (nullable = true)
Rpt Dist No: integer (nullable = true)	Part 1-2: integer (nullable = true)
Crm Cd: integer (nullable = true)	Crm Cd Desc: string (nullable = true)
Mocodes: string (nullable = true)	Vict Age: integer (nullable = true)
Vict Sex: string (nullable = true)	Vict Descent: string (nullable = true)
Premis Cd: integer (nullable = true)	Premis Desc: string (nullable = true)
Weapon Used Cd: integer (nullable = true)	Weapon Desc: string (nullable = true)
Status: string (nullable = true)	Status Desc: string (nullable = true)
Crm Cd 1: integer (nullable = true)	Crm Cd 2: integer (nullable = true)
Crm Cd 3: integer (nullable = true)	Crm Cd 4: integer (nullable = true)
LOCATION: string (nullable = true)	Cross Street: string (nullable = true)
LAT: double (nullable = true)	LON: double (nullable = true)

Για το Query 1 θα χρησιμοποίσουμε αρχικά το Dataframe API του Spark:

```
crime_data = Dataframe.withColumn("year", F.year("Date Rptd"))
crime_data = crime_data.withColumn("month", F.month("Date Rptd"))
count_df = crime_data.groupBy("year", "month").agg(F.count("*").alias("crime_count"))
window_spec = Window.partitionBy("year").orderBy(F.desc("crime_count"))
result = count_df.withColumn("ranking", F.row_number().over(window_spec)).filter("ranking <= 3") \
.orderBy("year", "ranking", "month")</pre>
```

Συγκεκριμένα, δημιουργούμε ένα νέο DataFrame από το αρχικό με δύο νέες στήλες month και year που δημιουργούμε με τις αντίστοιχες συναρτήσεις από την στήλη Date Rptd.

Στην συνέχεια με την μέθοδο groupBy αντιστοιχούμε κάθε ζεύγος year-month στον αντίστοιχο αριθμό εγκλημάτων που σημειώθηκαν και έπειτα δημιουργόυμε ένα βοηθητικό παράθυρο στο οπόιο θα δράσει η μέθοδος row_number για να πάρουμε το ranking κάθε μήνα ανά χρόνο.

Τέλος, φιλτράρουμε για να κρατήσουμε τους 3 μήνες κάθε χρόνου με τα περισσότερα εγκλήματα και εφμανίζουμε τα αποτελέσματα με την ζήτούμενη μορφή.

Την ίδια λογική θα χρησιμοποιήσουμε και με την χρήση του SQL API:

```
Dataframe.createOrReplaceTempView("Crime_data")

result = spark.sql("""

WITH RankedCrimeData AS (

SELECT

YEAR(`Date Rptd`) AS year,

MONTH(`Date Rptd`) AS month,

COUNT(*) AS `crime_total`,

ROW_NUMBER() OVER(PARTITION BY YEAR(`Date Rptd`) ORDER BY COUNT(*) DESC) as `#`

FROM Crime_data

GROUP BY year, month
)

SELECT year, month, crime_total, `#`
```

FROM RankedCrimeData

WHERE `#` <= 3

ORDER BY year, `#`, month;

111111

Στην παρακάτω φωτογραφία βλέπουμε και τα αποτελέσματα τα οποία είναι προφανώς τα ίδια για τις δύο υλοποιήσεις χωρίς να παρατηρούμε διαφορά στην επίδοση.

+	++		++
year	month	crime count	ranking
+	+		+
2010	3	17595	1
2010			
2010			
2011			
2011			
2011			
2012			
2012	: :		
2012	: :		
2013			
2013			
2013			
2014			
2014	: :		
2014			
2015			
2015			
2015		'	
2015			
2016			
	:		
2016			
2017			
2017			
2017			
2018			
2018			
2018			
2019			
2019			
2019			
2020			
2020			
2020			
2021			
2021			
2021			
2022			
2022			
2022			
2023			
2023			
2023	7	20216	3
+	+		++

Για το Query 2 αρχικά θα χρησιμοποίσουμε ξανά το Dataframe API του Spark:

```
street_crimes = Dataframe.filter(Dataframe['Premis Desc'] == 'STREET')
total_rows = street_crimes.count()
print(f"Total street crimes: {total_rows}")

street_crimes = street_crimes.withColumn(
    "Time of day",
    F.when((F.col("TIME OCC") >= 500) & (F.col("TIME OCC") < 1200), "morning")
    .when((F.col("TIME OCC") >= 1200) & (F.col("TIME OCC") < 1700), "noon")
    .when((F.col("TIME OCC") >= 1700) & (F.col("TIME OCC") < 2100), "afternoon")
    .when((F.col("TIME OCC") >= 2100) | (F.col("TIME OCC") < 500), "night")
    .otherwise("unknown")
)
street_crime_counts_by_time_of_day = street_crimes.groupBy("Time of day").agg(F.count("*").alias("crime_count")).orderBy("crime_count")

street_crime_counts_by_time_of_day.show()</pre>
```

Αρχικά φιλτράρουμε το Dataframe για να εξετάσουμε μόνο τα εγκλήματα που έλαβαν χώρα στον δρόμο. Έπειτα προσθέτουμε μια νέα στήλη με την ονομασία της χρονοπεριόδου που έγινε το έγκλημα με βάση την στήλη ΤΙΜΕ ΟCC. Τέλος, με την μέθοδο groupBy αντιστοιχόυμε κάθε περίοδο με τον αντίστοιχο αριθμό εγκλημάτων και εμφανίζουμε το αποτέλεσμα.

Την ίδια διαδικασία θα ακολοθήσουμε χρησιμοποιώντας το RDD API:

```
Dataframe_rdd = Dataframe.rdd

street_crimes_rdd = Dataframe_rdd.filter(lambda row: row['Premis Desc'] == 'STREET')

def time_of_day(hour):
    if 500 <= hour < 1200:
        return "morning"
```

```
elif 1200 <= hour < 1700:
    return "noon"
elif 1700 <= hour < 2100:
    return "afternoon"
elif 2100 <= hour or hour < 500:
    return "night"
else:
    return "unknown"

street_crimes_rdd_with_time_of_day = street_crimes_rdd.map(lambda row: row + ("Time of day", time_of_day(row['TIME OCC'])))

street_crime_counts_by_time_of_day_rdd = street_crimes_rdd_with_time_of_day.map(lambda row: (row[-</pre>
```

Συγκεκριμένα, αφόυ μετατρέψουμε το Dataframe σε κατάλληλη RDD μορφή, φιλτράρουμε πάλι κάθε row κρατώντας αυτά που αντιστοιχούν σε εγκλήματα που έγιναν στον δρόμο. Αφού ορίσουμε μια συνάρτηση που αντιστοιχίζει κατάλληλα τις τιμές της στήλης TIME OCC σε χρονοπεριόδους, εμπλουτίζουμε τις γραμμές με την νέα στήλη. Στην συνέχεια κάθε γραμμή αντιστοιχίζεται σε ένα ζευγάρι κλειδιού-τιμής, με κλειδί την χρονική περίοδο και η τιμή 1 και χρησιμοποιούμε την μέθοδο reduceByKey για μετρήσουμε τα αντίστοιχα εγκλήματα. Τέλος, ταξινομούμε το RDD σε αύξουσα σειρά.

1], 1)).reduceByKey(lambda x, y: x + y).sortBy(lambda x: x[1])

Τα αποτλέσματα που προκύπτουν από τις δύο υλοποιήσεις είναι τα ίδια με την RDD υλοποιήση να υστερεί σε ταχύτητα.

```
| Time of day|crime_count|
| Here of day|crime_count|
| Here of day|crime_count|
| Morning| 123128|
| Noon| 147303|
| Afternoon| 186208|
| Night| 236350|
```

```
[('morning', 123128),
('noon', 147303),
('afternoon', 186208),
('night', 236350)]
```

Για το Query 3 θα χρησιμοποίσουμε ξανά το Dataframe API του Spark:

```
crime_data_2015 = Dataframe.filter(F.year("Date Rptd") == 2015)
crime_data_2015 = crime_data_2015.filter(F.col("Vict Descent").isNotNull())
revgecoding_file = "./data/revgecoding.csv"
revgecoding_df = spark.read.csv(revgecoding_file, header=True, inferSchema=True)
Φιλτράρουμε το αρχικό Dataframe ώστε να κρατήσουμε τα εγκλήματα που έγιναν το 2015 και για
τα οποία υπάρχει καταγεγρανένη η καταγωγή του θύματος. Επίσης, δημιουργόυμε ένα νέο
Dataframe με βοηθητικές πληροφορίες για το zip code κάθε τοποθεσίας που μας ενδιαφέρει.
crime data_2015 = crime_data_2015.join(revgecoding_df.hint(join_hint), ['LAT', 'LON'],
'left_outer').select("Vict Descent", "ZIPcode")
crime_data_2015.show()
Στην συνέχεια πραγματοποιούμε join των δύο Dataframes με βάση την τοποθεσία και κρατάμε τις
στήλες που μας ενδιαφέρουν, δηλαδή την καταγωγή των θυμάτων και το zip code.
median income file = "./data/LA income 2015.csv"
median_income_df = spark.read.csv(median_income_file, header=True, inferSchema=True)
median_income_df = median_income_df.withColumn("Estimated Median Income",
F.regexp_replace(F.col("Estimated Median Income"), "[^\d.]", "").cast("int"))
top_incomes = median_income_df.orderBy(F.col("Estimated Median Income").desc()).limit(3)
lowest_incomes = median_income_df.orderBy(F.col("Estimated Median Income")).limit(3)
```

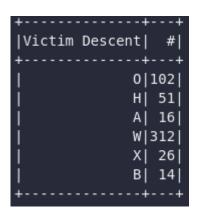
```
descents_low = crime_data_2015.join(lowest_incomes.hint(join_hint), F.col("ZIPcode") ==
F.col("Zip Code"), "inner")
descents_high = crime_data_2015.join(top_incomes.hint(join_hint), F.col("ZIPcode") == F.col("Zip Code"), "inner")
```

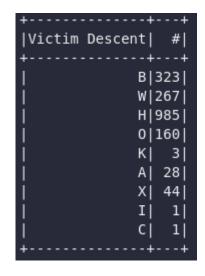
```
descents_low = descents_low.groupBy("Vict Descent").agg(F.count("*"))
descents_high = descents_high.groupBy("Vict Descent").agg(F.count("*"))
```

```
descents_low = descents_low.withColumnRenamed("Vict Descent", "Victim Descent")
descents_low = descents_low.withColumnRenamed("count(1)", "#")
```

```
descents_high = descents_high.withColumnRenamed("Vict Descent", "Victim Descent")
descents_high = descents_high.withColumnRenamed("count(1)", "#")
```

Αφού δημιουργήσουμε ένα νέο βοηθητικό Dataframe με πληροφορία για το μέσο εισόδημα και διαχωρίσουμε τις τρείς περιοχές με το υψηλότερο και χαμηλότερο εισόδημα κάνουμε join με βάση το zip code και αφου τα τροποποιήσουμε κατάλληλα εμφανίζουμε τα αποτελέσματα:





Για το Query 4 θα χρησιμοποίσουμε πάλι το Dataframe API:

Αρχικά φιλτράρουμε το Dataset κρατώντας μόνο τα εγκλήματα με την χρήση πυροβόλου όπλου.

```
gun_crimes = Dataframe.filter(F.col('Weapon Used Cd').like('1%'))
```

Αφου δημιουργήσουμε ένα Dataframe με τις κατάλληλες πληροφορίες για την τοποθεσία των διάφορων τμημάτων του Los Angeles κάνουμε join τα δύο Dataframes και προσθέτουμε μία νέα στήλη με την απόσατση του εγκλήματος από το αστυνομικό τμήμα που το ανέλαβε.

```
gun_crimes = gun_crimes.withColumn("year", F.year("Date Rptd"))
gun_crimes = gun_crimes.join(LAPD_stations_df.hint(join_hint), F.col("AREA") ==
F.col("PREC"), "inner").select("year", "DIVISION", "LAT", "LON", "X", "Y")
gun_crimes = gun_crimes.withColumn("Distance", get_distance_udf("LON", "LAT", "X",
"Y")).select("year", "DIVISION", "Distance")
```

Τέλος, ομαδοποιούμε ανά έτος και υπολογίζουμε το πλήθος και την μέση απόσταση αλλά και ανά αστυνομικό τμήμα.

```
result = gun_crimes.groupBy("year").agg(
    F.avg("Distance").alias("avg_distance"),
    F.count("Distance").alias("#")
).orderBy("year")
result.show()

result = gun_crimes.groupBy("DIVISION").agg(
    F.avg("Distance").alias("avg_distance"),
    F.count("Distance").alias("#")
).orderBy(F.col("#").desc())
result.show()
```

```
4.326341117149404
2011|2.7904269556134778|
2012 | 37.48698144373706 |
2013 | 2.830277105758418 |
2014 | 10.470155426256689 |
2015 2.7062554193223862
       2.71765650940967
2016
       4.33923298743777
2018 | 2.7355885132519364 |
       2.74082683654595
2019
                           7135
2020 | 8.613999633753169 |
2021 | 32.849532153520755 | 12101
2022 | 2.612071114474666 | 10067
2023 | 2.557452978934544 | 8484
```

```
DIVISION
                        avg distance
    77TH STREET | 5.745872330844248 | 16518
      SOUTHEAST |
                 13.56197627035242|13188|
         NEWTON| 9.896845692177651|
      SOUTHWEST | 4.160551099530668 |
     HOLLENBECK | 15.034085900271494 |
                                       6095
         HARBOR | 13.377359089703047 |
                                       5422
        RAMPART | 4.106286611009138 |
        MISSION| 7.545645264792472|
        OLYMPIC|1.8346273989145454|
                                       4305
      NORTHEAST | 10.459075190590411 |
       F00THILL | 3.810820653031428 |
      HOLLYWOOD | 12.119771738174023 |
        CENTRAL | 4.789187936164047 |
                                       3447
       WILSHIRE | 13.366513228684683 |
NORTH HOLLYWOOD | 14.03279776776525 |
                                       3339
    WEST VALLEY | 17.119572432381975 |
        PACIFIC | 13.276221563348491 |
                                       2640
       VAN NUYS| 2.215864249480073|
     DEVONSHIRE | 18.571732831383127 |
        TOPANGA | 3.4787633480390494 |
```

Για το δεύτερο σκέλος του ερωτήματος η στήλη Distance υπολογίζεται αυτήν την φορά με μια συνάρτηση που υπολογίζει την απόσταση από το πλησιέστερο τμήμα.

gun_crimes = gun_crimes.withColumn("Distance", find_closest_police_station_udf("LON",
"LAT")).select("year", "DIVISION", "Distance")

Παρακάτω βλέπουμε τα νέα αποτελέσματα:

```
Ivearl
            avq distance|
2010|3.9762380227744614|
                           8162
2011 | 2.4582677442240546 |
2012 37.13364209524676
                           6539
|2013|2.4592839268144866|
                           5851
2014 | 10.104120118078374 |
                           4882
2015 | 2.3883442547931346 |
                           6729
2016 | 2.425851163095081 |
                           8094
        4.00706930374014
                           7781
2018 | 2.411491332804097
2019 2.4303624630677043
                           7135
|2020| 8.304682016178997|
2021 32.53589290418105 12101
2022 | 2.3151258839871804 | 10067 |
2023 | 2.2715967638620214 |
```

```
DIVISION
                         avg distancel
     77TH STREET|
                    5.19309759370878 | 16518 |
       SOUTHEAST | 13.471602738284755 | 13188 |
          NEWTON| 9.444725167756307
       SOUTHWEST | 3.5601573646826825
                                        8615
      HOLLENBECK | 15.022619534245688 |
                                        6095 I
          HARBOR | 13.290425819008904 |
                                        5422
         RAMPART | 3.9680769025657825
                                        4976
         MISSION | 6.768267212875467
                                        4450
         OLYMPIC|1.7393567696353007
                                        4305
       NORTHEAST | 10.185143645838624 |
                                        3838
        F00THILL|3.7223003120431004|
                                        3761
       HOLLYWOOD | 12.112232296565534 |
                                        3541
         CENTRAL | 4.715568664043192
                                        3447
        WILSHIRE | 13.079051083443495 |
                                        3420
|NORTH HOLLYWOOD| 13.83065202357141|
                                        3339
    WEST VALLEY | 16.481518315328803 |
                                        2782
         PACIFIC | 13.190959733187475 |
                                        2640
        VAN NUYS | 2.1875146795048224 |
                                        2637
      DEVONSHIRE | 17.58786730446147
                                        2598
         TOPANGA | 3.173372840027675 |
```

Για τα joins των Queries 3, 4 χρησιμοποίησαμε joins hints αν και έχουμε μόνο 1 node. Παρόλα αυτά θεωρούμε πως το broadcast hint είναι το κατάλληλο για τα περισσότερα.

Στο BROADCAST join το μικρότερο Dataframe γίνεται broadcast σε όλους τους κόμβους και είναι ιδιαίτερα αποδοτικό όταν αυτό είναι αρκετά μικρό για να χωρέσει στην μνήμη των κόμβων.

Στην περίπτωσή μας τα joins πραγματοποιούνται με τα πολύ μικρότερα Dataframes LAPD_stations_df, top_incomes, lowest_incomes αλλά και σε τιμές (LAT, LON) που δεν είναι sorted για να χρησιμοποίησουμε κάποιο MERGE.

Ολόκληρος ο κώδικας υπάρχει στο https://github.com/VaggelisPap/AdvancedDB-Project