



# **ΤΜΗΜΑ ΠΛΗΡΟΦΟΡΙΚΗΣ**

---

# **DEPARTMENT OF INFORMATICS**

ΤΜΗΜΑ ΠΛΗΡΟΦΟΡΙΚΗΣ

ΜΕΡΟΣ Β ΕΡΓΑΣΙΑΣ Γ.Π.Σ

ΟΝ/ΝΥΜΟ: ΜΑΡΚΟ ΠΛΑΚΟΥ

Α.Μ: Π17107

ΗΜ/ΝΙΑ: 06/09/2021



# ΓΕΩΓΡΑΦΙΚΑ ΠΛΗΡΟΦΟΡΙΑΚΑ ΣΥΣΤΗΜΑΤΑ

# ΠΡΟΛΟΓΟΣ



Στην εργασία  
αυτή:

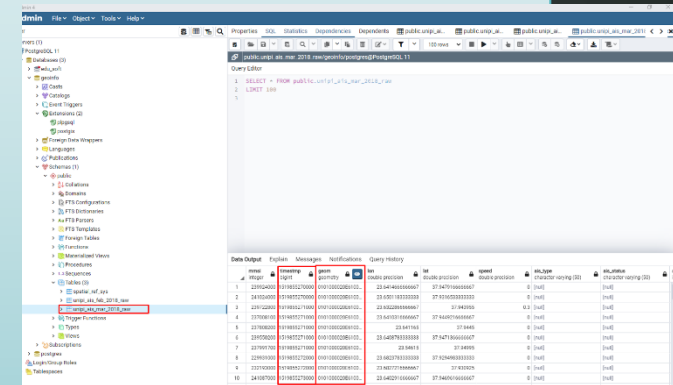
Παρουσιάζονται τα  
μέρη 1,2,3 της  
εργασίας.

Επεξηγείται η  
απάντηση κάθε  
ερωτήματος.

Χρησιμοποιήθηκαν τα  
εργαλεία  
Postgres, Qgis, Python3  
(Jupyter NB)

# 1. Φόρτωση δεδομένων (data loading)

- ▶ Χρησιμοποιήθηκε η βάση δεδομένων **Saronic Gulf Maritime AIS Dataset**.
- ▶ Η εισαγωγή της στήλης **point** έγινε μέσω της Postgres με **SIRD 4326**. Τα σήματα αποθηκεύτηκαν σε διαφορετικούς πίνακες για κάθε μήνα.

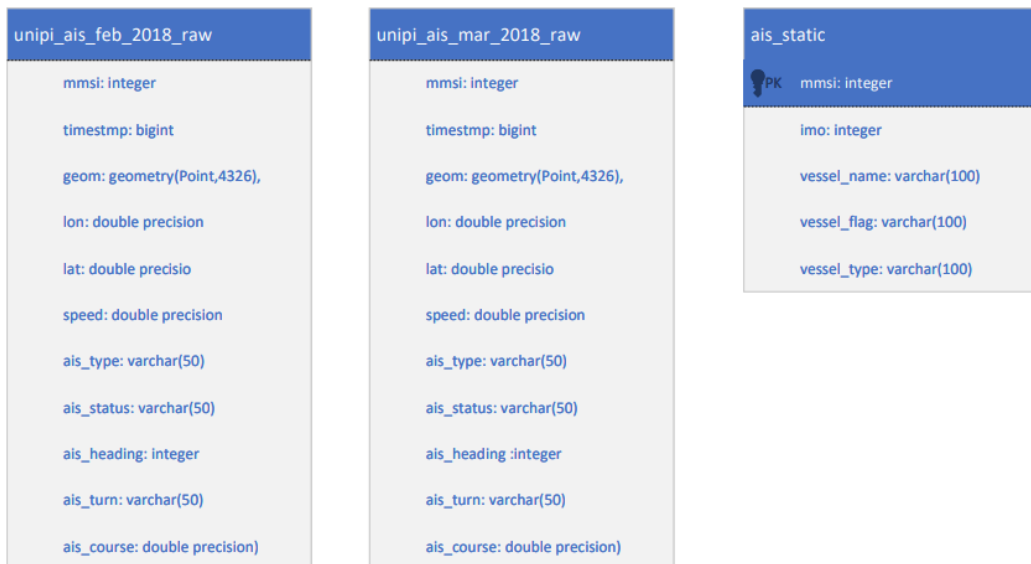


```
mar_geodf.head()
```

	mmsi	timestamp	geom	lon	lat	speed	ais_type	ais_status	ais_heading	ais_turn	ais_course
239569	636013190	1519855200000	POINT (23.53934 37.88567)	23.539338	37.885675	1.4	None	None	44.0	None	15.5
239571	237991700	1519855201000	POINT (23.54615 37.94995)	23.546150	37.949950	0.0	None	None	NaN	None	268.0
239572	239722800	1519855201000	POINT (23.63220 37.94397)	23.632205	37.943968	0.0	None	None	NaN	None	0.0
239573	239550200	1519855201000	POINT (23.64085 37.94716)	23.640850	37.947157	0.0	None	None	264.0	None	142.1
239570	241024000	1519855201000	POINT (23.65010 37.93166)	23.650097	37.931657	0.0	None	None	332.0	None	178.1

# 1. ΣΧΕΔΙΑΓΡΑΜΜΑ ΤΗΣ ΒΑΣΗΣ

## PART B: B1.GEOINFO UML DATABASE SCHEMA



# 2.Γνωριμία με τα δεδομένα και προετοιμασία (data acquaintance and preprocessing)

## 2.1 ΚΑΘΑΡΙΣΜΟΣ ΤΩΝ ΔΕΔΟΜΕΝΩΝ

- ▶ Αφαιρέθηκαν οι διπλότυπες εγγραφές, και αυτές που είχαν NaN τιμές στις στήλες `timestamp`, `speed`.
- ▶ Αφαιρέθηκε ο χωρικός θόρυβος.
- ▶ Έμειναν μόνο οι εγγραφές που το `mmsi` τους υπήρχε στα στατικά δεδομένα και αυτές των οποίων η ταχύτητα είναι μικρότερη των 50 κόμβων.
- ▶ Υπολογίστηκε η επιτάχυνση.



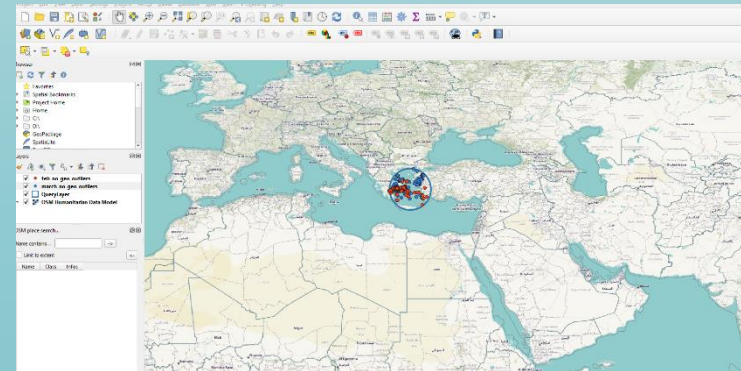
## ΜΕΤΑ ΤΟΝ ΚΑΘΑΡΙΣΜΟ

Το πλήθος δεδομένων του Φεβρουαρίου μειώθηκε κατά **1.145.481** εγγραφές.

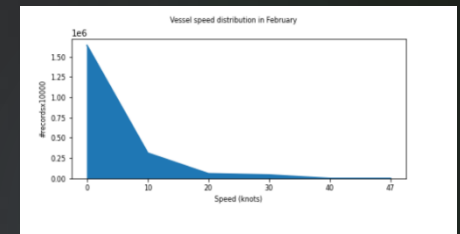
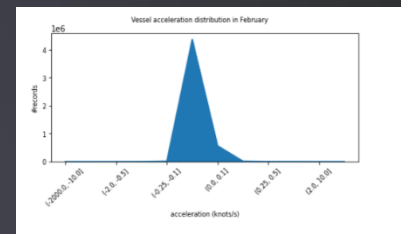
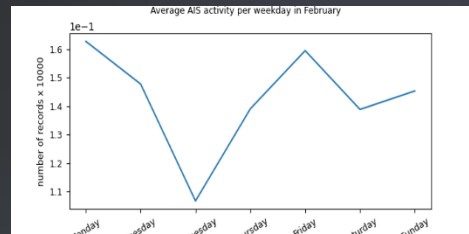
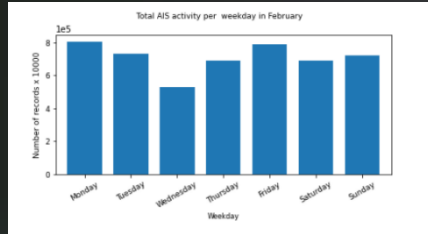
```
len(final_feb_geodf) #the final length of the February dataset.  
4960765
```

Το πλήθος του Μαρτίου μειώθηκε κατά **1.587.841**.

```
len(final_mar_geodf) #the final length of the March dataset.  
5116283
```



► **2.2-2.3** Τα διαγράμματα δειγματοληψίας και ταχύτητας για τον μήνα Φεβρουάριο:



► **2.4** Η ύπαρξη του R-tree ευρετηρίου μείωσε κατά πολύ τον χρόνο εκτέλεσης :

PRIN(421147.105 ms)

META(613.951 ms)

```
Query Editor Query History

1 explain analyze
2 select *
3 from (select box2d(geon) as box2d
4       from unip1_atls_feb_2018_clean
5       where ut_content(instr(transform(st_geomfromtext('polygon((23.423558 37.866290, 23.473977 37.8365)
6 unip1_atls_feb_2018_clean
7 where foo.box2d && geom);

Data Output Explain Messages Notifications Geometry Viewer

QUERY PLAN
1 text
2 Nested Loop (cost=1000.00:2476768.11 rows=16533 width=311) (actual=1,
3 Join Filter: (box2d(unip1_atls_feb_2018_clean_1.geon) & geometry & unip1
4 Rows Removed by Join Filter: 69467724
5 => Seq Scan on unip1_atls_feb_2018_clean (cost=0.00:91266.78 rows=132,
6 Maintenance (cost=1000.00:236028.50 rows=39 width=42) (actual time=
7 = Gather (cost=0.00:0.00:236028.14 rows=38 width=32) (actual time=2.
8 Workers Planned: 2
9 Workers Launched: 2
10 => Parallel Seq Scan on unip1_atls_feb_2018_clean unip1_atls_feb_201
11 Filter: (GistIndex <= '1517434001000_30jg' and GistIndex <= '1
12 Rows Removed by Filter: 1653541
13 Planning Time: 0.861 ms
14 Execution Time: 42147.135 ms
```

```

1 explain analyze
2 select +
3 from (select box2d(geom) as box2d
4       from unip1_at15_feb_2018_clean
5       where st_contains(st_transform(st_geomfromtext('polygon((23.423538 37.868298, 23.472977 37.8656
6         unip1_at15_feb_2018_clean
7       where foo_box2d && geom);

```

Data Output Explain Messages Notifications Geometry Viewer

---

**QUERY PLAN**

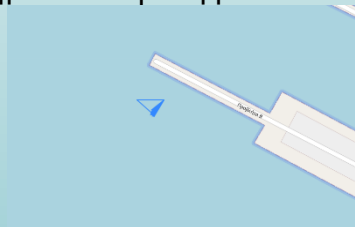
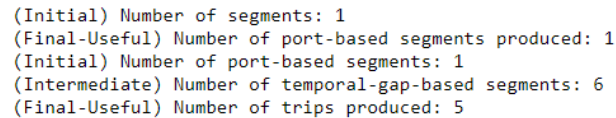
test

- 1 Gather (cost=1020.15..799841.95 rows=704429 width=311) (actual time=304.284..613.888 rows=...
- 2 Workers Planned: 2
- 3 Workers Launched: 2
- 4 → Parallel Loop (cost=20.15..72899.05 rows=203812 width=511) (actual time=386.302..507.16...
- 5 → Parallel Seq Scan on unip\_at15\_feb\_2018\_clean unip\_at15\_feb\_2018\_clean\_1 (cost=0.00..64...
- 6 Filter (BtreeScan ~ 151743601000 bkgnd) And (BtreeScan ~ 1517645201000 bkgnd)
- 7 Rows Removed by Filter: 1632541
- 8 → Bitmap Heap Scan on unip\_at15\_feb\_2018\_clean (cost=20.15..1343.85 rows=496 width=44...
- 9 Recheck Cond: (box2d(unip\_at15\_feb\_2018\_clean\_1.geom)) geometry && geom)
- 10 Heap Blocks: exact=76
- 11 → Bitmap Index Scan on idx\_unip\_at15\_feb\_2018 (cost=0.00..20.02 rows=496 width=0) (ac...
- 12 Index Cond: (box2d(unip\_at15\_feb\_2018\_clean\_1.geom)) geometry && geom)
- 13 Planning Time: 305.913 ms



10

- ### ▶ 3.2 Διάφορα χωροχρονικά ερωτήματα πραγματοποιήθηκαν:



```

1 select *
2 from traktartoken_bowl;
3 insert into traktartoken_bowl(traktartoken) as traktartoken_bowl on row
4 values(traktartoken_bowl) as traktartoken_bowl on row
5 where traktartoken_bowl.traktartoken = traktartoken_bowl.traktartoken;

```

- ### 3.3 Temporal alignment and Resampling Using Helper

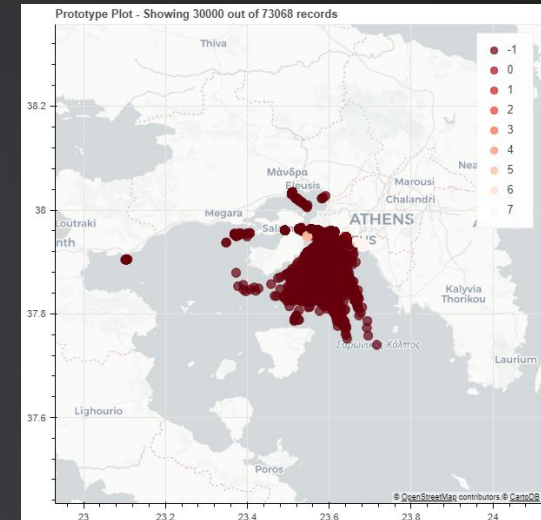


# 3.4-3.5

## 3.4 Υλοποιήθηκε ο OPTICS και βρήκε 9 συστάδες:

```
In [14]: gdf_radians_resampled_and_aligned = np.radians(ra_df[['lat', 'lon']]) #creating the radians
optics_resampled_and_aligned = OPTICS(max_eps=1/6371, min_samples=len(ra_df)//50, metric='haversine').fit(gdf_radians_resampled_and_aligned)
set(optics_resampled_and_aligned.labels_)
helper2.get_clusters_centers(gdf_radians_resampled_and_aligned, optics_resampled_and_aligned.labels_)

Out[14]: array([[37.94691256, 23.64038312],
 [37.94731144, 23.63740902],
 [37.9427454, 23.64046713],
 [37.94420766, 23.63208578],
 [37.95457064, 23.55113343],
 [37.94993253, 23.54616344],
 [37.9316622, 23.65010351],
 [37.93275919, 23.6812189 ]])
```



## 3.5 Υλοποιήθηκε ο Evolving Clusters και βρήκε 14 evolving cluster αποτελούμενα από 1 πλοίο:

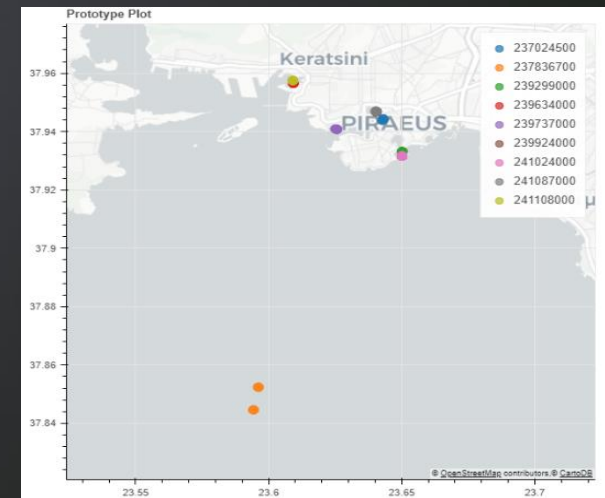
### Implementing Evolving Clusters Algorithm

```
In [12]: [res_mcs, res_mc] = evolving_clusters(final, distance_threshold=1000, min_cardinality=1, time_threshold=1, disable_progress_bar=False)

100% [██████████] 71910/71910 [06:03<00:00, 197.94it/s]

In [13]: for row in res_mcs.iteruples(): # display the clusters found
          print(row)

Pandas(Index=0, clusters=(241024000), st=Timestamp('2018-02-01 01:07:00'), et=Timestamp('2018-02-01 01:08:00'))
Pandas(Index=0, clusters=(241024000), st=Timestamp('2018-02-06 01:23:00'), et=Timestamp('2018-02-06 01:24:00'))
Pandas(Index=0, clusters=(241087000), st=Timestamp('2018-02-14 18:37:00'), et=Timestamp('2018-02-14 18:38:00'))
Pandas(Index=0, clusters=(241024000), st=Timestamp('2018-02-17 00:24:00'), et=Timestamp('2018-02-17 00:25:00'))
Pandas(Index=0, clusters=(241024000), st=Timestamp('2018-02-17 18:06:00'), et=Timestamp('2018-02-17 18:07:00'))
Pandas(Index=0, clusters=(237836700), st=Timestamp('2018-02-18 13:57:00'), et=Timestamp('2018-02-18 13:58:00'))
Pandas(Index=0, clusters=(239299000), st=Timestamp('2018-03-03 14:13:00'), et=Timestamp('2018-03-03 14:14:00'))
Pandas(Index=0, clusters=(241024000), st=Timestamp('2018-03-06 18:26:00'), et=Timestamp('2018-03-06 18:27:00'))
Pandas(Index=0, clusters=(239634000), st=Timestamp('2018-03-13 00:34:00'), et=Timestamp('2018-03-13 00:36:00'))
Pandas(Index=0, clusters=(241108000), st=Timestamp('2018-03-15 07:56:00'), et=Timestamp('2018-03-15 07:57:00'))
Pandas(Index=0, clusters=(239737000), st=Timestamp('2018-03-19 22:08:00'), et=Timestamp('2018-03-19 22:09:00'))
Pandas(Index=0, clusters=(239737000), st=Timestamp('2018-03-20 01:34:00'), et=Timestamp('2018-03-20 01:35:00'))
Pandas(Index=0, clusters=(237024500), st=Timestamp('2018-03-28 09:52:00'), et=Timestamp('2018-03-28 09:53:00'))
Pandas(Index=0, clusters=(239924000), st=Timestamp('2018-03-29 03:54:00'), et=Timestamp('2018-03-29 03:56:00'))
```





# Data Story και Κατακλείδα

Η εργασία αυτή ήταν εξαιρετικά ενδιαφέρουσα αφού έθεσε πραγματικά ερωτήματα και επίσης χρησιμοποιήθηκαν αληθινά δεδομένα.

Τα **συμπεράσματα** που εξήχθησαν και 'διηγούνται' την **ιστορία των δεδομένων** είναι:

- ▶ Μειωμένη κίνηση των πλοίων τις Τετάρτες του Φεβρουαρίου και τις Κυριακές του Μαρτίου.
- ▶ Δεν υπήρχαν πλοία με ακριβώς κοινή τροχιά.
- ▶ Από την εύρεση των hot-spots, είναι εμφανές ότι η πλειοψηφία των σημάτων λαμβάνεται σε κοντινή απόσταση από το λιμάνι του Πειραιά.

## Εν κατακλείδι.

- ▶ Έχοντας ως 'πρώτη ύλη' τα ακατέργαστα δεδομένα με θόρυβο που δεν παρέχουν σχεδόν καμία χρήσιμη πληροφορία και γνώση, με την κατάλληλο 'καθαρισμό' και προ επεξεργασία και με την χρήση τεχνικών αναλυτικής δεδομένων καταλήγουμε σε χρήσιμη πληροφορία που μπορεί να αξιοποιηθεί κατάλληλα από τους ειδικούς (domain experts).