

ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ

ΠΜΣ «ΠΛΗΡΟΦΟΡΙΑΚΑ ΣΥΣΤΗΜΑΤΑ ΚΑΙ ΥΠΗΡΕΣΙΕΣ»

ΚΑΤΕΥΘΥΝΣΗ «ΜΕΓΑΛΑ ΔΕΔΟΜΕΝΑ ΚΑΙ ΑΝΑΛΥΤΙΚΗ»

Διαχείριση Δεδομένων για Σχεσιακές και μη Σχεσιακές Βάσεις Δεδομένων



ΕΥΑΓΓΕΛΙΔΑΚΗΣ ΛΕΑΝΔΡΟΣ

Περιεχόμενα

Περιγραφή του προβλήματος	2
Περιγραφή των δεδομένων	2
Μέθοδοι και τεχνικές που χρησιμοποιήθηκαν	4
Εισαγωγή δεδομένων στη MongoDB	5
Ευρετήρια	7
Ανάλυση των δεδομένων	8
Αλιευτική δραστηριότητα	8
Συσταδοποίηση τροχιών αλιευτικών πλοίων	13
Συσταδοποίηση σημείων αλιευτική δραστηριότητας	18
Παράνομη αλιεία	22
Συμπεράσματα	24

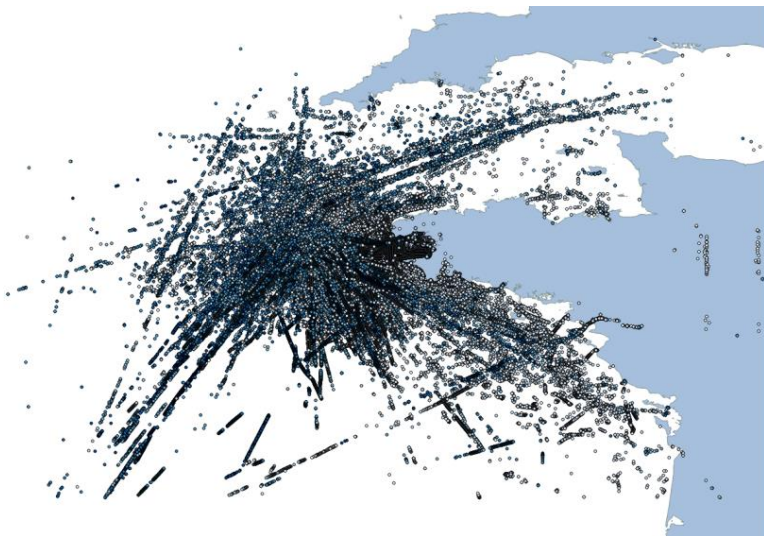
Περιγραφή του προβλήματος

Στην παρούσα εργασία εξετάζουμε τα διαθέσιμα δεδομένα από AIS (Automatic Identification System), για το εξάμηνο από 1^η Οκτωβρίου 2015 έως τον Μάρτιο του 2016, που αναφέρονται στην περιοχή της Κελτικής Θάλασσας και του Βискаϊκού κόλπου. Το επίκεντρο της μελέτης μας είναι η αλιεία. Εξετάζουμε τα αλιευτικά πλοία και κάνουμε συσταδοποίηση ανάλογα με την τροχιά που κινούνται (trajectory clustering) καθώς και σημειακή συσταδοποίηση για να προσδιορίσουμε πιθανές περιοχές ψαρέματος (fishing spots clustering). Τέλος, εντοπίζουμε περιπτώσεις παράνομης αλιείας (fishing violations) σε περιοχές όπου απαγορεύεται.

Περιγραφή των δεδομένων

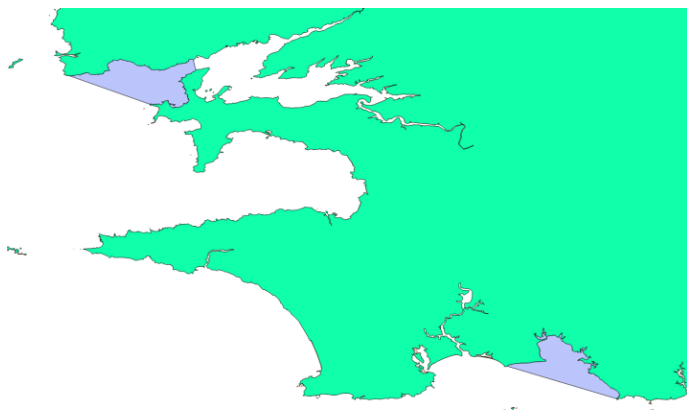
Το σύνολο των δεδομένων περιέχει δεδομένα πλοήγησης, δεδομένα σχετικά με τα πλοία, γεωγραφικά δεδομένα αλλά και περιβαλλοντολογικά. Τα δεδομένα που εξετάζουμε κυρίως, είναι το σύνολο [P1] AIS Data, που περιέχει μεταξύ άλλων το αρχείο `nari_dynamic.csv` και το `nari_static.csv`. Εκεί περιέχονται οι περισσότερες καταχωρήσεις των σημάτων μέσω του AIS, όπως λαμβάνονται από τους σταθμούς που τα αποθηκεύουν. Το σύστημα AIS χρησιμοποιεί 27 είδη μηνυμάτων (ITU 1, ITU 2 κ.λπ.)

Στο αρχείο `nari_dynamic.csv` υπάρχουν πληροφορίες που στέλνουν αυτόματα τα πλοία κατά τη διάρκεια της δραστηριότητάς τους, και αναφέρονται στην τοποθεσία τους, την κατεύθυνση που κινούνται, την ταχύτητα κλπ. Αυτά αντιστοιχούν στα μηνύματα τύπου ITU 1, 2, 3, 18 και 19. Ο χρόνος (timestamp) που αναφέρεται σε κάθε μήνυμα αντιστοιχεί στην στιγμή που ο σταθμός τα λαμβάνει και όχι τη στιγμή που τα εκπέμπει το κάθε πλοίο. Επίσης κάθε πλοίο αναφέρει με έναν ακέραιο (navigationalstatus) την κατάστασή του τη στιγμή εκείνη (π.χ. εν κινήσει, αγκυροβολημένο κ.λπ.). Σε αυτή τη μεταβλητή, το 7 αντιστοιχεί στην αλιευτική δραστηριότητα και έτσι μπορούμε και εξάγουμε εύκολα τα δεδομένα που αφορούν την αλιεία.



Εικόνα 1. Τα δεδομένα του nari_dynamic.csv

Στο αρχείο Fishing constraints, περιέχονται δύο περιοχές που απαγορεύεται η αλιεία, στην περιοχή της Βρετάνης. Η πρώτη περιοχή (Camaret, βόρεια), αναφέρεται στο διάστημα από 05/02/2015 έως 21/01/2016, ενώ η δεύτερη (Concarneau, νότια) από 03/09/2015 έως 15/10/2015.



Εικόνα 2. Περιοχές όπου απαγορεύεται η αλιεία

Μέθοδοι και τεχνικές που χρησιμοποιήθηκαν

Η αποθήκευση των δεδομένων γίνεται χρησιμοποιώντας τη μη-σχεσιακή βάση δεδομένων MongoDB, σε συνδυασμό με την Python για την εισαγωγή, εξαγωγή και ανάλυσή τους.

Για την αποθήκευση, επεξεργαζόμαστε ελάχιστα τα δεδομένα του `pari_dynamic.csv` και τα εισάγουμε στη MongoDB. Αυτό γίνεται χρησιμοποιώντας κυρίως τις βιβλιοθήκες `pandas` και `pymongo` της Python. Την ίδια διαδικασία ακολουθούμε και για τα δεδομένα για τις περιοχές με απαγόρευση αλιείας. Και στις δύο περιπτώσεις, κάνουμε χρήση του γεγονότος ότι η MongoDB δέχεται έγγραφα τύπου GeoJSON, δηλαδή γεωγραφικά δεδομένα όπως συντεταγμένες (Points), πολύγωνα (Polygons) κ.λπ. και γι'αυτό εισάγουμε τα δεδομένα αφότου τα προσαρμόσουμε στη μορφή αυτή.

Καθότι μελετάμε την αλιευτική δραστηριότητα και τον έλεγχο αυτής σε συγκεκριμένες περιοχές (χρονικά εξαρτώμενες), δημιουργούμε ευρετήρια με σκοπό να βελτιστοποιήσουμε την ταχύτητα των διαδικασιών. Χρησιμοποιούμε απλά ευρετήρια (single indexes) αλλά και τα “compound” ευρετήρια της MongoDB, δηλαδή ευρετήρια με πολλαπλές παραμέτρους. Συγκεκριμένα, φτιάχνουμε compound ευρετήριο ως προς το χρόνο και ως προς την τοποθεσία (2d sphere).

Για την ανάλυση, γίνεται προσπάθεια η ανάκτηση των απαιτούμενων δεδομένων να προέρχεται κυρίως από queries στη βάση, καθώς ο αρχικός σκοπός της εργασίας είναι η ενασχόληση με μη-σχεσιακές βάσεις δεδομένων. Έτσι με την Python αντλούμε τα δεδομένα από τη Mongo, με τρόπο τέτοιο ώστε να είναι έτοιμα για την μετέπειτα ανάλυσή τους. Η ανάλυση περιλαμβάνει την παρουσίαση κάποιων βασικών πληροφοριών που εξάγουμε σχετικά με την αλιευτική δραστηριότητα, την συσταδοποίηση των τροχιών των αλιευτικών πλοίων (fishing vessels trajectory clustering), την συσταδοποίηση των μεμονωμένων σημείων αλιείας (fishing spots clustering) και τέλος τον εντοπισμό παράνομης αλιευτικής δραστηριότητας (fishing violations).

Η συσταδοποίηση γίνεται χρησιμοποιώντας αλγόριθμους βελτιστοποιημένους για Big Data. Συγκεκριμένα, χρησιμοποιούμε τους αλγορίθμους [KMeans](#), [DBSCAN](#), [OPTICS](#) και [BIRCH](#). Οι DBSCAN και OPTICS είναι μέθοδοι με βάση την πυκνότητα των δεδομένων (density based) και χρησιμοποιούνται ευρέως κατά την ανάλυση γεωχωρικών (geospatial) δεδομένων.

Ο BIRCH χρησιμοποιεί δενδρικές δομές για τη συσταδοποίηση και εφαρμόζεται γιατί παρέχει εντυπωσιακές ταχύτητες σε μεγάλα δεδομένα (πρόκειται για βελτίωση του MiniBatch-KMeans).

Για την συσταδοποίηση σημείων, οι διαδικασίες είναι απλές καθώς πρόκειται για κλασσικού τύπου συσταδοποίηση. Η συσταδοποίηση τροχιών, διαφέρει αρκετά από την προηγούμενη διαδικασία. Η τροχιά που σχηματίζει ένα πλοίο αποτελείται από τις προηγούμενες τοποθεσίες του σε χρονολογική σειρά. Τα σχήματα που δημιουργούν οι τροχιές είναι αρκετά πολύπλοκα ιδιαίτερα για τα αλιευτικά πλοία. Επίσης, όπως είναι και λογικό, οι τροχιές των αλιευτικών πλοίων διαφέρουν κατά πολύ συγκριτικά με τα υπόλοιπα πλοία. Για να συγκρίνουμε δύο τροχιές, πρέπει να βρούμε έναν τρόπο μέτρησης της ομοιότητας δύο τροχιών, δηλαδή δύο σχημάτων. Ένας από τους τρόπους να το κάνουμε αυτό είναι η απόσταση του [Hausdorff](#) που μετράει την απόσταση μεταξύ δύο υποσυνόλων ενός μετρικού χώρου εξετάζοντας τις αποστάσεις μεταξύ των επιμέρους στοιχείων τους. Έτσι χρησιμοποιούμε αυτή την απόσταση ως είσοδο στους αλγορίθμους συσταδοποίησης, κατασκευάζοντας χειροκίνητα τον πίνακα των αποστάσεων.

Για τον εντοπισμό παράνομης αλιευτικής δραστηριότητας, κάνουμε γεωχωρικά queries για συγκεκριμένα χρονικά διαστήματα, κάνοντας χρήση της δυνατότητας της MongoDB να εντοπίζει σημεία που βρίσκονται σε κάποιο γεωμετρικό σχήμα (π.χ. Polygons). Με την υποστήριξη και από το compound ευρετήριο που δημιουργούμε, τέτοιου είδους queries εκτελούνται αποδοτικά.

Εισαγωγή δεδομένων στη MongoDB

Αρχικά δημιουργούμε μια βάση δεδομένων 'AISdata' στη MongoDB. Η εισαγωγή των επιμέρους collections γίνεται από την Python μέσω της βιβλιοθήκης pymongo. Για την εισαγωγή των δεδομένων, διαβάζουμε τα δεδομένα του nari_dynamic ως dataframe μέσω της βιβλιοθήκης pandas. Παρακάτω φαίνεται ο αντίστοιχος κώδικας και οι πέντε πρώτες γραμμές των δεδομένων.

```

1 import pymongo
2 from pymongo import MongoClient
3 import pandas as pd
4 import seaborn as sns
5 import geopandas as geo
6 import pprint

```

```

1 client = MongoClient()
2 client = MongoClient('localhost', 27017)
3 db = client.AISdata
4 nari_dynamic = db.nari_dynamic

```

```

1 data = pd.read_csv("../data/[P1] AIS data/nari_dynamic.csv")
2 data.head()

```

	sourceemsi	navigationalstatus	rateofturn	speedoverground	courseoverground	trueheading	lon	lat	t
0	245257000	0.0	0.0	0.1	13.1	36	-4.465718	48.382490	1443650402
1	227705102	15.0	-127.0	0.0	262.7	511	-4.496571	48.382420	1443650403
2	228131600	15.0	-127.0	8.5	263.7	511	-4.644325	48.092247	1443650404
3	228051000	0.0	-127.0	0.0	295.0	511	-4.485108	48.381320	1443650405
4	227574020	15.0	-127.0	0.1	248.6	511	-4.495441	48.383660	1443650406

Όπως βλέπουμε ο χρόνος είναι σε μορφή ακεραίου (UNIX epochs). Καθώς είναι δύσκολη η ανθρώπινη κατανόηση σε αυτή τη μορφή τον μετατρέπουμε στη συνήθη μορφή. Επίσης, θα ενώσουμε το γεωγραφικό μήκος και πλάτος ώστε να το φέρουμε στη μορφή του GeoJSON που αναγνωρίζει η MongoDB. Η μορφή αυτή είναι τύπου: <field>: {type: <GeoJSON type>, coordinates: <coordinates>}. Για παράδειγμα: { type: "Point", coordinates: [40, 5] }. Έπειτα διαβάζουμε γραμμή προς γραμμή το αρχείο και το εισάγουμε στη βάση.

```

1 data['t']=pd.to_datetime(data.t,unit='s')

```

```

1
2 def add_to_db():
3
4     for index, row in data.iterrows():
5         point=[row.lon,row.lat]
6         doc = {'sourceemsi':row.sourceemsi,
7               'navigationalstatus':row.navigationalstatus,
8               'rateofturn':row.rateofturn,
9               'speedoverground':row.speedoverground,
10              'courseoverground':row.courseoverground,
11              'trueheading':row.trueheading,
12              't':row.t,
13              'location':{'type':'Point',
14                          'coordinates':point}}
15         nari_dynamic.insert_one(doc)
16
17     return 'Done'
18
19

```

```

1 add_to_db()

```

```
'Done'
```

Για την εισαγωγή των απαγορευμένων περιοχών εξετάζουμε τα δεδομένα που βρίσκονται στο αρχείο fishing_interdiction.shp. Για να διαβάσουμε το αρχείο χρησιμοποιούμε τη βιβλιοθήκη geopandas. Σε αυτό βρίσκονται δύο καταχωρήσεις, για τις δύο περιοχές αντίστοιχα. Η πρώτη (Camaret) είναι ένα απλό πολύγωνο (Polygon) ενώ η δεύτερη (Concarneau) είναι πολύ-

πολύγωνο (MultiPolygon). Για την εισαγωγή ενός πολυγώνου ως GeoJSON γίνεται με τον ίδιο τρόπο που είδαμε προηγουμένως. Η μόνη διαφορά είναι ότι εδώ πρέπει να αναπροσαρμόσουμε τον τρόπο που διαβάζει η geopandas τα δεδομένα, ώστε αντί για tuples να χρησιμοποιεί λίστες, ώστε να μπορούν να εισαχθούν απευθείας στη βάση. Για την περίπτωση του MultiPolygon, παίρνουμε ένα-ένα τα πολύγωνα που περιέχει και επαναλαμβάνουμε την προηγούμενη διαδικασία.

```
1 forbidden_fishing_areas=geo.read_file("../data/[C5] Fishing Constraints/fishing_interdiction.shp")
2
3 #Camaret (Polygon)
4 x0,y0=forbidden_fishing_areas.geometry[0].exterior.coords.xy
5 polygon0=list(list(t) for t in zip(x0,y0))
6
7 #Concarneau (MultiPolygon)
8 polygons1=[]
9 multipolygon=forbidden_fishing_areas.geometry[1]
10 for i in range(len(multipolygon)):
11     x1,y1=multipolygon[i].exterior.coords.xy
12     polygon1=list(list(t) for t in zip(x1,y1))
13     polygons1.append(polygon1)
14
15 #Insert forbidden areas to mongodb :
16
17 forbidden_fishing_areas = db.forbidden_fishing_areas
18
19 #Camaret
20 forbidden_area_0 = {'id':'Camaret',
21                    'type':'Polygon',
22                    'coordinates':[polygon0]}
23 forbidden_fishing_areas.insert_one(forbidden_area_0)
24
25 #Concarneau
26 i=0
27 for poly in polygons1:
28     forbidden_area_1 = {'id':'Concarneau'+str(i),
29                        'type':'Polygon',
30                        'coordinates':[poly]}
31     forbidden_fishing_areas.insert_one(forbidden_area_1)
32     i+=1
```

Ευρετήρια

Όπως αναφέραμε, θα χρησιμοποιήσουμε δύο ευρετήρια. Ένα compound index που περιλαμβάνει πρώτα το χρόνο ("t") και έπειτα το γεωγραφικό σημείο ("location"), και ένα απλό ευρετήριο ως προς την κατάσταση που δηλώνει κάθε πλοίο ("navigationstatus").

Για την δημιουργία τους εκτελούμε απευθείας στη βάση της ακόλουθες δύο εντολές.

```
db.nari_dynamic.createIndex({"t":1,"location" : "2dsphere"})
```

```
db.nari_dynamic.createIndex({"navigationstatus":1})
```


Ανάλυση των δεδομένων

Αλιευτική δραστηριότητα

Αρχικά θα εξάγουμε κάποια βασικά στοιχεία σχετικά με την αλιευτική δραστηριότητα από τα διαθέσιμα δεδομένα του αρχείου `nari_dynamic.csv`. Τα queries που χρησιμοποιούμε γίνονται εφαρμόζοντας το **aggregation pipeline** της MongoDB, που αποτελείται από διάφορα στάδια.

Πλήθος καταχωρήσεων για αλιευτική δραστηριότητα για κάθε πλοίο

```
1 #Find distinct fishing vessels and count their fishing activity instances
2
3 query = ([
4     { '$match': {'navigationalstatus': {'$eq': '7'}}},
5     { '$group': {'_id': '$sourcemsi', 'num_appearances': {'$sum': 1}}},
6     { '$sort': {'num_appearances': -1}}
7 ])
8
9 document = nari_dynamic.aggregate(query)
10
11 results=[]
12 for i in document:
13     results.append(i)
14
15
16 res=pd.DataFrame(results).copy()
17 res['_id']=res['_id'].astype(str)
```

```
1 fishing_vessels=list(res['_id'].astype(int))
2 print('There are',len(fishing_vessels),'vessels that were fishing')
```

There are 208 vessels that were fishing

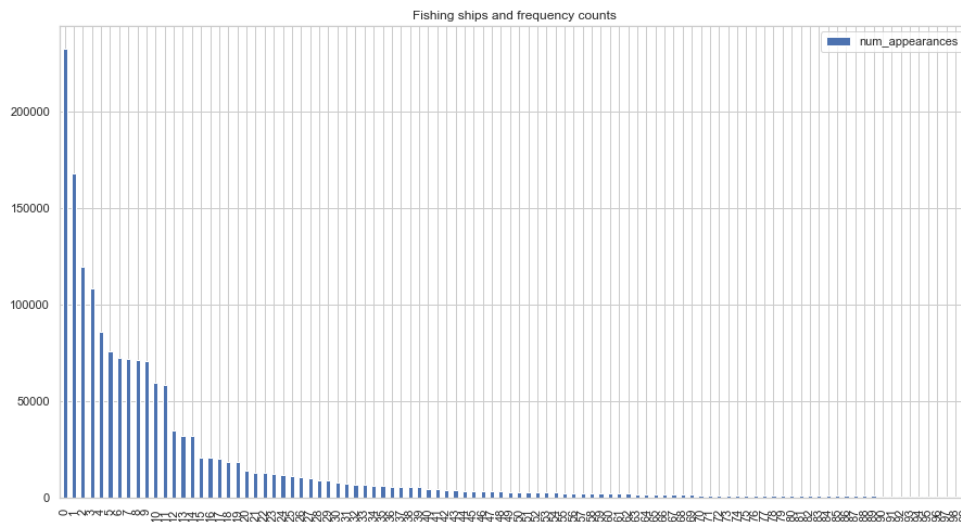
Όπως φαίνεται, στέλνουμε το query στη βάση και αποθηκεύουμε σε dataframe τα δεδομένα που επιστρέφονται. Επίσης εμφανίζουμε το πλήθος των διακριτών πλοίων που δήλωσαν ότι ψαρεύανε, το οποίο είναι 208 πλοία. Παρακάτω βλέπουμε τις πέντε πρώτες και πέντε τελευταίες γραμμές των αποτελεσμάτων αλλά και ένα ραβδόγραμμα με τα 100 πρώτα ποτελέσματα.

```
1 res
```

	_id	num_appearances
0	227941000	232591
1	227300000	167941
2	228394000	119457
3	227391000	108516
4	228144000	85724
...
203	235103397	2
204	250002464	2
205	225361000	1
206	228143600	1
207	227573960	1

```
1 res[:100].plot(kind='bar',figsize=(15,8),title='Fishing ships and frequency counts')
2
```

<matplotlib.axes._subplots.AxesSubplot at 0x242849a2668>



Παρατηρούμε ότι τα πρώτα δέκα πλοία έχουν περισσότερες από 50.000 καταχωρήσεις. Αντίθετα υπάρχουν πλοία με ελάχιστες (έως και 1). Στη συνέχεια παραθέτουμε κάποια στατιστικά σχετικά με το πλήθος των καταχωρήσεων.

```
1 res.describe()
```

	num_appearances
count	208.0000
mean	8071.0577
std	26218.5440
min	1.0000
25%	33.0000
50%	453.5000
75%	2735.7500
max	232591.0000

Όπως φαίνεται, κατά μέσο όρο τα αλιευτικά πλοία έχουν περίπου 80.000 καταχωρήσεις με τυπική απόκλιση περίπου 26.000. Η ελάχιστη τιμή είναι 1 καταχώρηση ενώ η μέγιστη πάνω από 230.000. Επίσης τα μισά πλοία έχουν το πολύ 453 καταχωρήσεις ενώ το 75% αυτών το πολύ 2.736

Χώρες προέλευσης

Το αναγνωριστικό (source mmsi) κάθε πλοίου αντιστοιχεί σε έναν 9ψήφιο ακέραιο αριθμό. Τα πρώτα τρία στοιχεία αυτού του αριθμού αντιστοιχούν στη χώρα προέλευσης. Έτσι μπορούμε να υπολογίσουμε πόσα αλιευτικά αντιστοιχούν σε κάθε κωδικό, δηλαδή σε κάθε χώρα.

```
: 1 mmsi=res['_id'].apply( lambda x : ''.join(list(x)[0:3]) )
  2 print('There are vessels from',len(mmsi.unique()),'countries that were fishing')
  3 print('Country codes: ',mmsi.unique())

There are vessels from 11 countries that were fishing
Country codes:  ['227' '228' '226' '224' '219' '232' '235' '244' '218' '225' '250']
```

Οι χώρες που αντιστοιχούν στους παραπάνω κωδικούς είναι :

227, 228, 226: Γαλλία

224, 225 : Ισπανία

218: Γερμανία

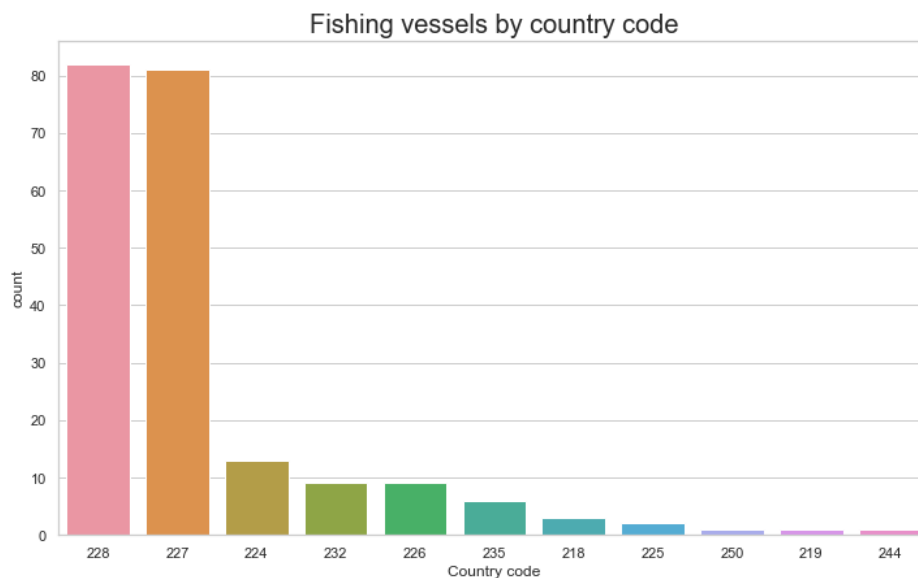
219: Δανία

232, 235 : Ηνωμένο Βασίλειο

244: Ολλανδία

250: Ιρλανδία

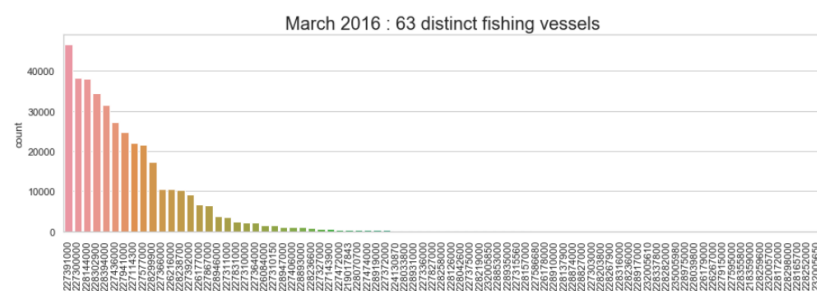
```
1 plt.figure(figsize=(12,7))
2 plt.title('Fishing vessels by country code',fontsize=20)
3 sns.countplot(mmsi,order = mmsi.value_counts().index)
4 plt.xlabel('Country code')
5 plt.show()
```

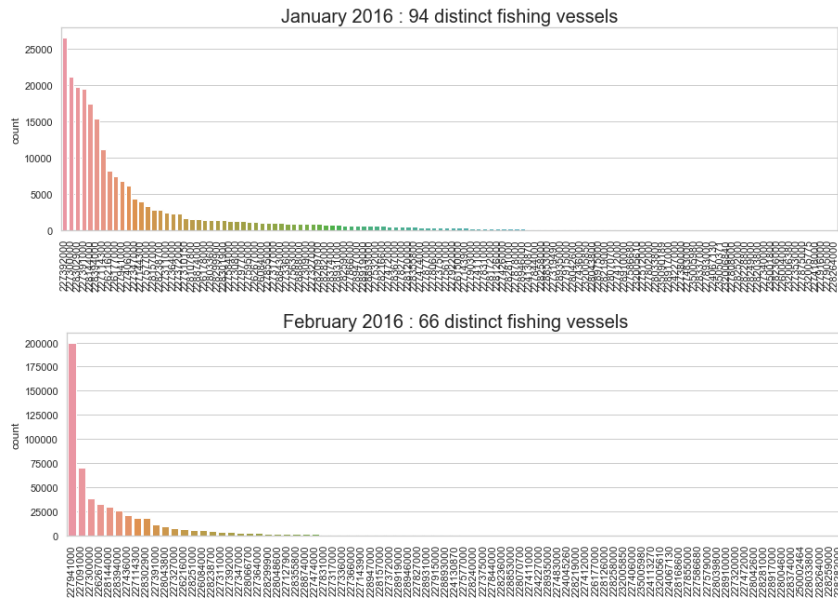


Στοιχεία ανά μήνα

Το χρονικό εύρος των δεδομένων είναι από τον Σεπτέμβρη του 2015 έως τον Μάρτιο του 2016. Στα επόμενα, υπολογίζουμε πόσα πλοία ενεπλάκησαν σε αλιευτική δραστηριότητα για κάθε μήνα, και φτιάχνουμε ένα διάγραμμα για κάθε έναν. Κάνοντας χρήση της μεθόδου **\$month** της MongoDB, μπορούμε εύκολα να ομαδοποιήσουμε τα έγγραφα με βάση τους μήνες και με τη μέθοδο **\$push**, αποθηκεύουμε σε ένα διάνυσμα όλα τα αναγνωριστικά των πλοίων που εμφανίστηκαν τον κάθε μήνα.

```
1 #Get all coordinates from the fishing vessels,group by months and count
2
3 query = ([
4     {'$match':{'navigationalstatus':{'$eq':7}}
5     },
6     {'$group': {'_id':{'month':{'$month':'$t'}}},
7       'count':{'$sum':1},
8       'sourceemsi':{'$push':{'$sourceemsi'}}
9     },
10    {'$sort':{'_id':1}}
11  ])
12
13 document = nari_dynamic.aggregate(query,allowDiskUse=True)
14
15 results=[]
16 for i in document:
17     results.append(i)
18
19 res_monthly=pd.DataFrame(results)
20
21 months={
22     1:'January 2016',
23     2:'February 2016',
24     3:'March 2016',
25     9:'September 2015',
26     10:'October 2015',
27     11:'November 2015',
28     12:'December 2015',
29 }
30 i=0
31
32 for month in res_monthly['_id']:
33
34     monthly_counts=pd.Series(res_monthly.sourceemsi[i]).value_counts()
35
36     plt.figure(figsize=(12,5))
37     plt.title(months[month['month']]+' : '+str(len(np.unique(monthly_counts)))+ ' distinct fishing vessels',fontsize=12)
38     plot=sns.countplot(res_monthly.sourceemsi[i],order=monthly_counts.index)
39     plot.set_xticklabels(plot.get_xticklabels(),rotation=90)
40     plt.show()
41     i+=1
```





Συσταδοποίηση τροχιών αλιευτικών πλοίων

Αρχικά, πρέπει να φέρουμε τα δεδομένα στη μορφή που χρειάζεται για να μπορούμε να τα αναλύσουμε. Αυτή είναι τύπου: #πλοίο : [[ΘΕΣΗ 1],[ΘΕΣΗ2],...].

Δηλαδή πρέπει να αντιστοιχήσουμε σε κάθε διακριτό πλοίο μια λίστα με τις προηγούμενες συντεταγμένες του. Τα πλοία που θα χρησιμοποιήσουμε τα έχουμε βρει στα προηγούμενα. Από αυτά, επιλέγουμε να αφαιρέσουμε τα πλοία με λίγες εγγραφές (<20) καθώς και αυτά με υπερβολικά πολλές (>50.000). Εν τέλει, θα αναλύσουμε τις τροχιές για 154 πλοία.

Το query που χρησιμοποιούμε χρησιμοποιεί και εδώ το aggregation pipeline. Με τη χρήση του τελεστή \$in επιλέγουμε τα πλοία που βρίσκονται στη λίστα με το σύνολο των πλοίων και με τη χρήση του \$push, εισάγουμε σε μια λίστα όλες τις τοποθεσίες για κάθε πλοίο.

```

1 #Get all coordinates from the fishing vessels
2
3 query = ([
4     {'$match':{'sourceemsi':{'$in':fishing_vessels[-164:-10:1]}},
5     {'$group': {'_id':'$sourceemsi', 'locations':{'$push':'$location.coordinates'}}},
6     {'$sort':{'_id':-1}}
7 ])
8
9 document2 = nari_dynamic.aggregate(query,allowDiskUse=True)
10
11 results2=[]
12 for i in document2:
13     results2.append(i)
14
15
16 res2=pd.DataFrame(results2)
17 print(res2)
18 res2['_id']=res2['_id'].astype(str)

```

```

      id                      locations
0  226150000  [[-2.8469417, 46.808846], [-2.8375932999999995...
1  228120000  [[-2.610005, 45.885033], [-2.6056633, 45.88513...
2  227303000  [[-5.3681116, 48.111694], [-5.3650199999999995...
3  228240000  [[-4.841038, 47.8540340000000006], [-4.81325, 4...
4  228203800  [[-5.135842, 48.7179150000000005], [-5.13125, 4...
..
149 228236000  [[-5.07866140000000005, 47.633907], [-5.0735765...
150 228021600  [[-4.92382140000000005, 48.291733], [-4.918665,...
151 228106800  [[-4.6752234, 47.962368], [-4.6761484000000001,...
152 226004000  [[-2.02976660000000003, 46.298626], [-1.923525,...
153 232005800  [[-5.9387784, 49.2931], [-5.9356465, 49.29647]...

```

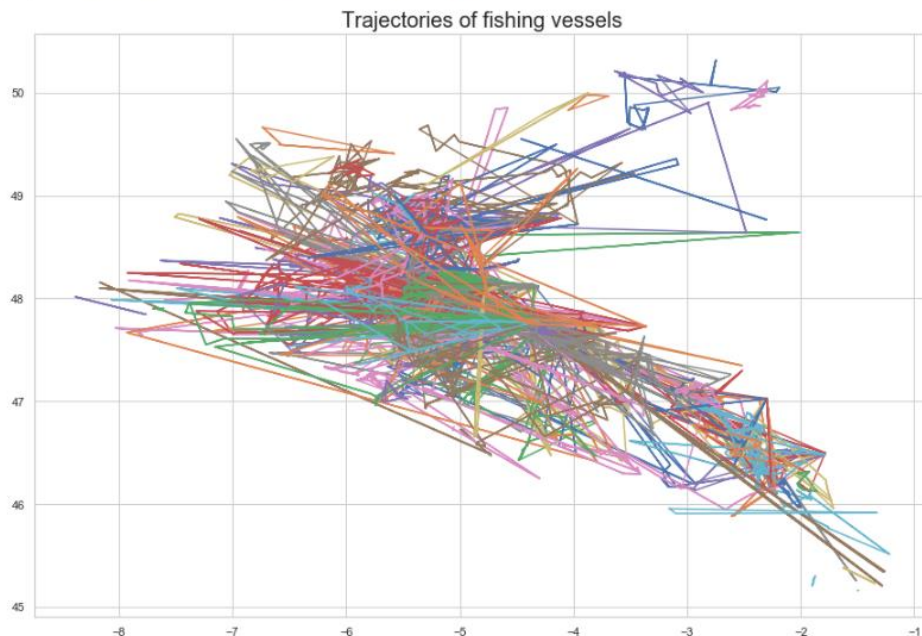
[154 rows x 2 columns]

Παρακάτω φαίνονται οι επιμέρους τροχιές που σχηματίζονται από τα δεδομένα.

```

1 plt.figure(figsize=(15,10))
2 plt.title('Trajectories of fishing vessels',fontsize=20)
3 for trajectory in res2.locations:
4     plt.plot(trajectory[:, 0], trajectory[:, 1])

```



Τώρα θα χρησιμοποιήσουμε την απόσταση Hausdorff για να υπολογίσουμε τον πίνακα αποστάσεων, στον οποίο μετά θα εφαρμόσουμε τους αλγορίθμους συσταδοποίησης. Η απόσταση αυτή υπάρχει [διαθέσιμη](#) από τη βιβλιοθήκη SciPy και έτσι η διαδικασία γίνεται αρκετά εύκολα. Παρακάτω φαίνεται ο κώδικας για τη δημιουργία του πίνακα αποστάσεων και ο κώδικας για μια συνάρτηση όπου θα οπτικοποιεί την ανάθεση των συστάδων. Επειδή κάποιος αλγόριθμος (DBSCAN), σχηματίζουν μία συστάδα με τα outliers (ή noise), στην οποία μάλιστα ανατίθεται ο αριθμός -1 , επιλέγουμε να τη χρωματίσουμε με μαύρο χρώμα ώστε να ξεχωρίζει.

```
1 #Creating hausdorff distance matrix
2
3 def hausdorff( u, v):
4     d = max(directed_hausdorff(u, v)[0], directed_hausdorff(v, u)[0])
5     return d
6
7 traj_count = len(res2.locations)
8 D = np.zeros((traj_count, traj_count))
9
10
11 for i in range(traj_count):
12     for j in range(i + 1, traj_count):
13         distance = hausdorff(res2.locations[i], res2.locations[j])
14         D[i, j] = distance
15         D[j, i] = distance
16
17
18 #Function for plotting clustered trajectories
19
20 colors = plt.rcParams['axes.prop_cycle'].by_key()['color']
21 colors.extend(['firebrick', 'olive', 'indigo', 'khaki', 'teal', 'saddlebrown',
22               'skyblue', 'coral', 'darkorange', 'lime', 'darkorchid', 'dimgray'])
23
24 def plot_cluster(trajs, clusters):
25
26     cluster_count = np.max(clusters) + 1
27
28     for traj, cluster in zip(trajs, clusters):
29
30         if cluster == -1:
31             plt.plot(traj[:, 0], traj[:, 1], c='k', linestyle='dashed')
32
33         else:
34             plt.plot(traj[:, 0], traj[:, 1], c = colors[cluster % len(colors)])
35
36     plt.show()
```

Για να πάρουμε μια πρώτη γνώμη σχετικά με το βέλτιστο πλήθος συστάδων, χρησιμοποιούμε τη μέθοδο του αγκώνα (elbow method) με τον αλγόριθμο KMeans, για τον πίνακα αποστάσεων που υπολογίσαμε προηγουμένως.


```

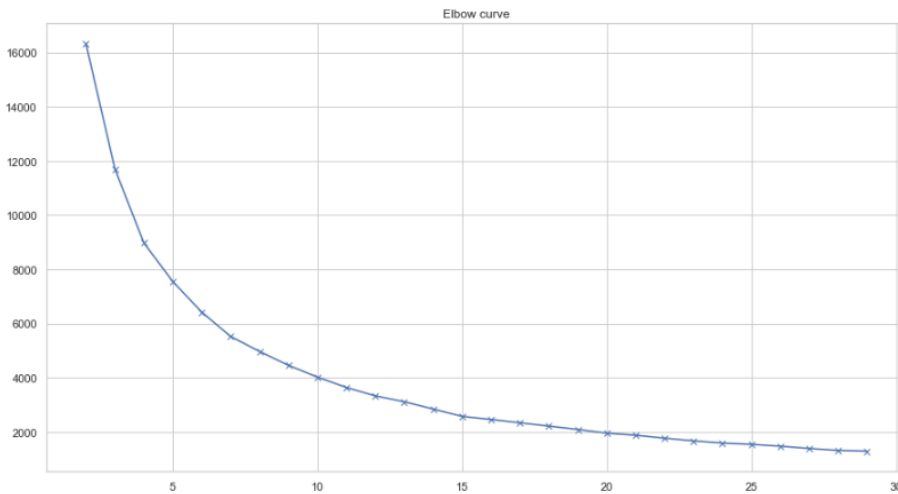
1 from sklearn.cluster import KMeans
2
3 def elbow_method(data,max_clusters=10):
4     inertias = []
5     for k in range(2, max_clusters):
6         kmeans = KMeans(n_clusters=k,random_state=7,n_jobs=-1)
7         kmeans.fit(data)
8         inertias.append(kmeans.inertia_)
9
10    fig = plt.figure(figsize=(15, 8))
11    plt.plot(range(2, max_clusters), inertias, 'bx-')
12    plt.grid(True)
13    plt.title('Elbow curve')

```

```

1 elbow_method(D,30)

```

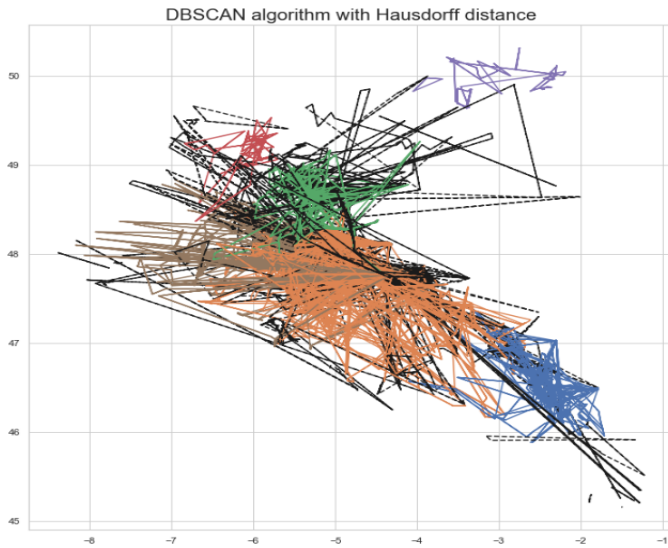


Όπως φαίνεται, το βέλτιστο πλήθος σύμφωνα με το παραπάνω διάγραμμα φαίνεται να είναι από 6 έως 10 συστάδες. Στη συνέχεια εφαρμόζουμε τους αλγορίθμους συσταδοποίησης. Οι αλγόριθμοι DBSCAN και OPTICS δεν απαιτούν την εισαγωγή του πλήθους των συστάδων. Παρόλ'αυτά, ο πρώτος καταλήγει σε 7 συστάδες ενώ ο δεύτερος σε 10. Κάτι που φαίνεται να έρχεται σε συμφωνία με τα αποτελέσματα της μεθόδου του αγκώνα. Στον αλγόριθμο BIRCH, επιλέγουμε 7 συστάδες.

DBSCAN

```
1 dbSCAN = DBSCAN(eps=4, min_samples=4,n_jobs=-1)
2 clusters = dbSCAN.fit_predict(D)
3 print('Nof clusters:',len(np.unique(clusters)))
4 plt.figure(figsize=(12,12))
5 plt.title('DBSCAN algorithm with Hausdorff distance',fontsize=20)
6 plot_cluster(res2.locations, clusters)
7
```

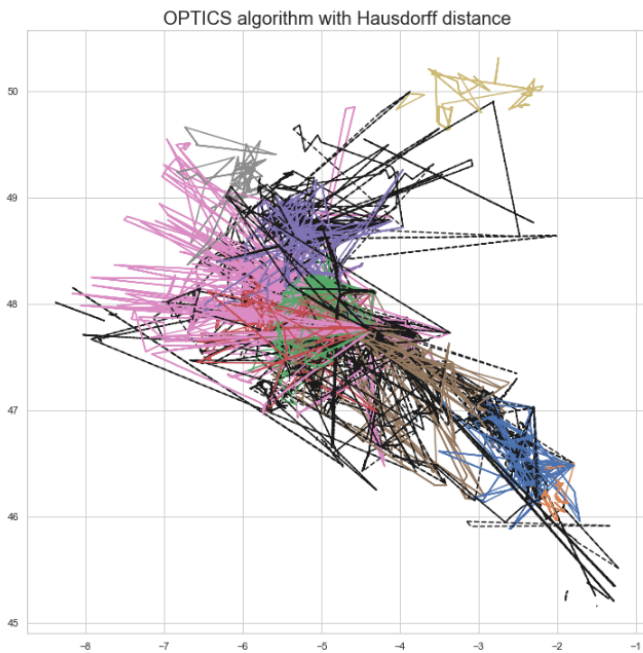
Nof clusters: 7



OPTICS

```
1 optics = OPTICS(eps=5, min_samples=6,leaf_size=30,n_jobs=-1)
2 clusters = optics.fit_predict(D)
3 print('Nof clusters:',len(np.unique(clusters)))
4 plt.figure(figsize=(12,12))
5 plt.title('OPTICS algorithm with Hausdorff distance',fontsize=20)
6 plot_cluster(res2.locations, clusters)
7
```

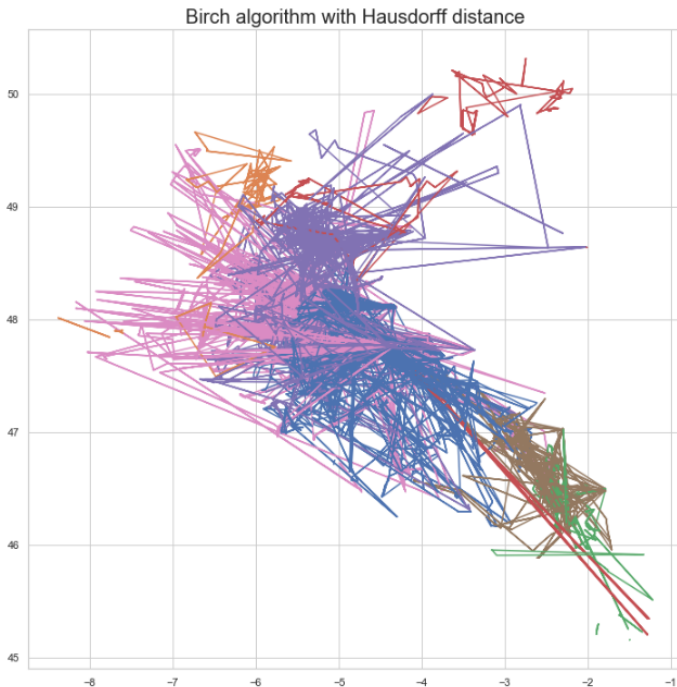
Nof clusters: 10



BIRCH

```
1 brc = Birch(n_clusters=7)
2 clusters = brc.fit_predict(D)
3 print('Nof clusters:', len(np.unique(clusters)))
4 plt.figure(figsize=(12,12))
5 plt.title('Birch algorithm with Hausdorff distance', fontsize=20)
6 plot_cluster(res2.locations, clusters)
7
```

Nof clusters: 7



Συσταδοποίηση σημείων αλιευτική δραστηριότητας

Στη συνέχεια εξετάζουμε τη συσταδοποίηση των σημείων που εκδηλώθηκε αλιευτική δραστηριότητα. Εδώ δεν θα υπολογίσουμε για τα αλιευτικά πλοία όλες τις προηγούμενες γεωγραφικές τους θέσεις, αλλά μόνο τις θέσεις στις οποίες ενεπλάκησαν σε αλιευτικές διαδικασίες.

Και πάλι χρησιμοποιούμε το aggregation pipeline όπως κάναμε και στην προηγούμενη ενότητα, για να πάρουμε τις τοποθεσίες από 154 πλοία. Παρακάτω φαίνεται ο αντίστοιχος κώδικας καθώς και ένα διάγραμμα με τις τοποθεσίες. Εδώ δεν μας ενδιαφέρει το κάθε πλοίο, αλλά συνολικά οι συντεταγμένες, γι'αυτό τις ενώνουμε σε μία ενιαία λίστα ώστε να τις δούμε όλες μαζί.

```

1 #Get all coordinates from the fishing vessels
2
3 query = {
4     '$match': {'navigationstatus': {'$eq': 7},
5               'sourcemmsi': {'$in': fishing_vessels[-164:-10:1]}},
6     '$group': {'_id': '$sourcemmsi',
7               'locations': {'$push': '$location.coordinates'}}},
8 }
9
10 document = nari_dynamic.aggregate(query, allowDiskUse=True)
11
12 results = []
13 for i in document:
14     results.append(i)
15
16
17 res3 = pd.DataFrame(results)
18 print(res3)
19 res3['_id'] = res3['_id'].astype(str)

```

```

      _id      locations
0  227472000  [[-5.22182, 48.238859999999995], [-5.0097933, ...
1  228840000  [[-2.2973516, 47.02714], [-2.2978283999999998, ...
2  227317000  [[-4.8319782999999999, 47.862957], [-4.85794, 4...
3  228282000  [[-5.6109519999999999, 47.398586], [-5.517648, ...
4  228209700  [[-2.3717099999999998, 46.34022], [-2.3722167, ...
...
149 228259000  [[-3.0895984, 47.11665], [-3.0897317, 47.11651...
150 228917000  [[-1.2608683, 45.346214], [-1.2799083, 45.3450...
151 232007350  [[-2.94427, 49.90136], [-2.94427, 49.90136]]
152 228252000  [[-4.1755949999999995, 47.835873], [-4.175598, ...
153 227315560  [[-5.183143, 48.78235], [-5.3549120000000001, 4...

```

[154 rows x 2 columns]

```

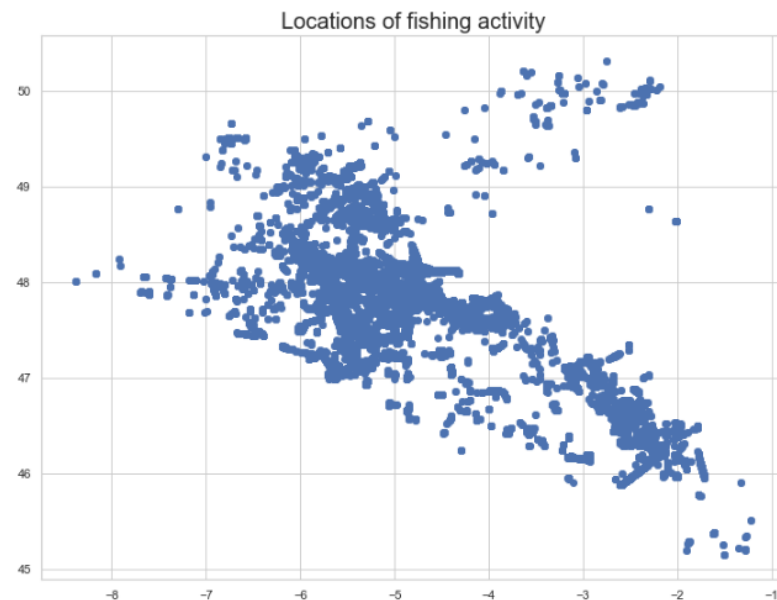
1 fishing_coords = []
2 for sublist in res3.locations:
3     for item in sublist:
4         fishing_coords.append([item[0], item[1]])
5
6 fishing_coords = pd.DataFrame.from_records(fishing_coords, columns=['lon', 'lat'])

```

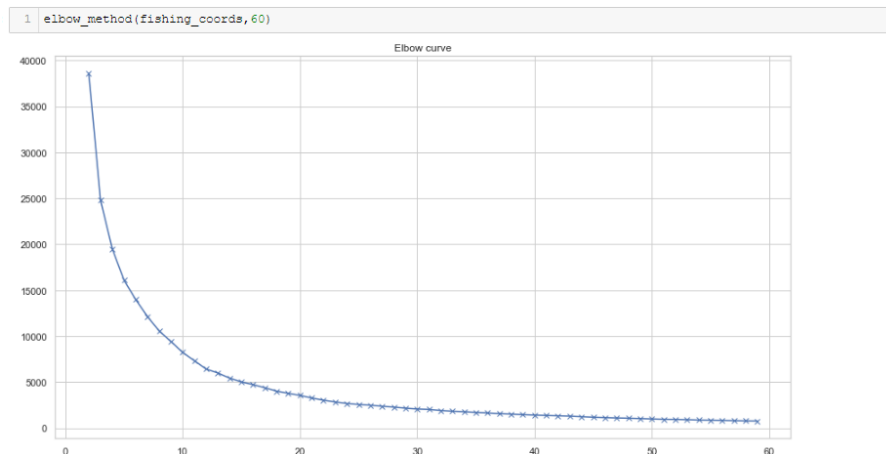
```

1 plt.figure(figsize=(12,9))
2 plt.title('Locations of fishing activity', fontsize=20)
3 plt.scatter(fishing_coords.lon, fishing_coords.lat)
4 plt.show()

```



Στη συνέχεια φαίνονται τα αποτελέσματα του elbow method για τα δεδομένα. Και πάλι το βέλτιστο πλήθος συστάδων φαίνεται να είναι κοντά στο 10.

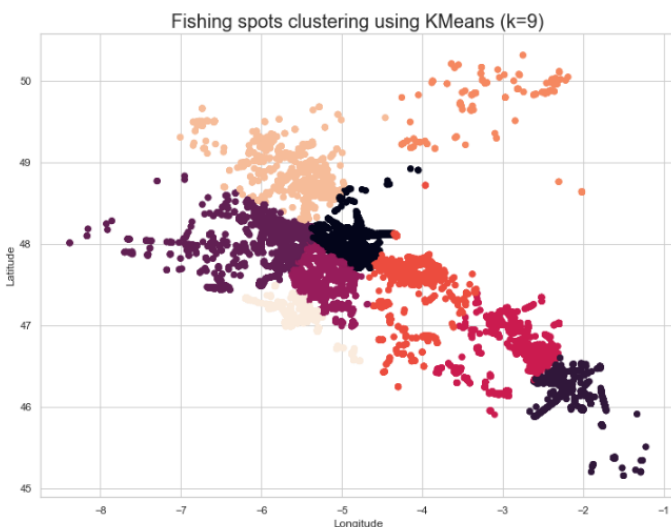


Έπειτα χρησιμοποιούμε τους αλγόριθμους KMeans και Birch για τη συσταδοποίηση. Οι άλλοι αλγόριθμοι DBSCAN και OPTICS δεν μας έδωσαν ικανοποιητικά αποτελέσματα και οι χρόνοι εκτέλεσής τους ήταν ιδιαίτερα μεγάλοι. Γι'αυτό δεν τους χρησιμοποιούμε σε αυτή την περίπτωση.

KMEANS

```
1 kmeans = KMeans(9,n_jobs=-1)
2 clusters = kmeans.fit_predict(fishing_coords)
3 print('Nof clusters:',len(np.unique(clusters)))
4
5 clustered_fishing_coords=fishing_coords.copy()
6 clustered_fishing_coords['clusters']=clusters
7
8 fig, ax = plt.subplots(figsize=[12, 9])
9 plt.scatter(fishing_coords.lon,fishing_coords.lat,c=clusters)
10 plt.title('Fishing spots clustering using KMeans (k=9)',fontsize=20)
11 ax.set_xlabel('Longitude')
12 ax.set_ylabel('Latitude')
13 plt.show()
```

Nof clusters: 9

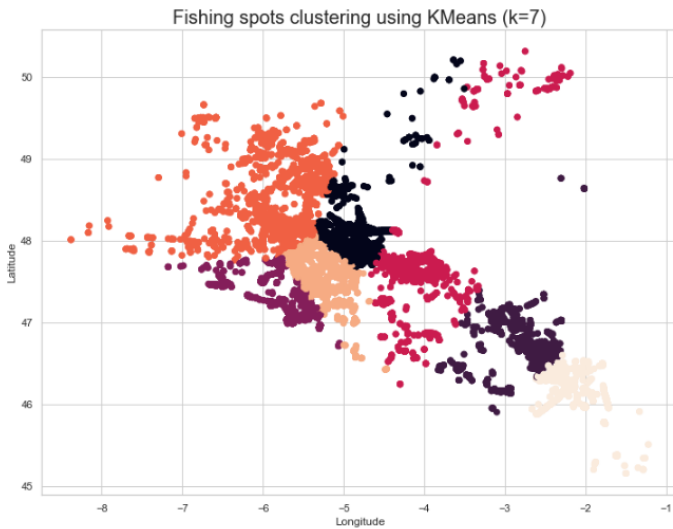


```

1 kmeans = KMeans(7,n_jobs=-1)
2 clusters = kmeans.fit_predict(fishing_coors)
3 print('Nof clusters:',len(np.unique(clusters)))
4
5 clusterd_fishing_coors=fishing_coors.copy()
6 clusterd_fishing_coors['clusters']=clusters
7
8 fig, ax = plt.subplots(figsize=[12, 9])
9 plt.scatter(fishing_coors.lon,fishing_coors.lat,c=clusters)
10 plt.title('Fishing Spots clustering using KMeans (k=7)',fontsize=20)
11 ax.set_xlabel('Longitude')
12 ax.set_ylabel('Latitude')
13 plt.show()

```

Nof clusters: 7



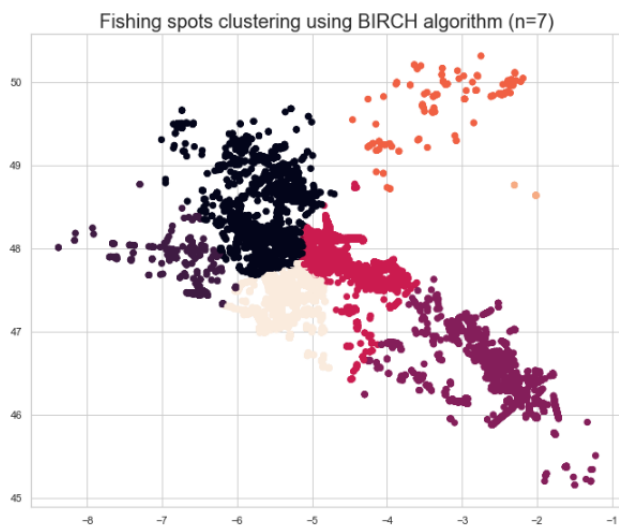
BIRCH

```

1 #Using BIRCH
2
3 brc = Birch(n_clusters=7)
4 clusters = brc.fit_predict(fishing_coors)
5 print('Nof clusters:',len(np.unique(clusters)))
6 clusterd_fishing_coors=fishing_coors.copy()
7 clusterd_fishing_coors['clusters']=clusters
8 plt.figure(figsize=(11,9))
9 plt.title('Fishing spots clustering using BIRCH algorithm (n=7)',fontsize=20)
10 plt.scatter(clusterd_fishing_coors.lon,clusterd_fishing_coors.lat,c=clusterd_fishing_coors['clusters'])
11 plt.show()

```

Nof clusters: 7

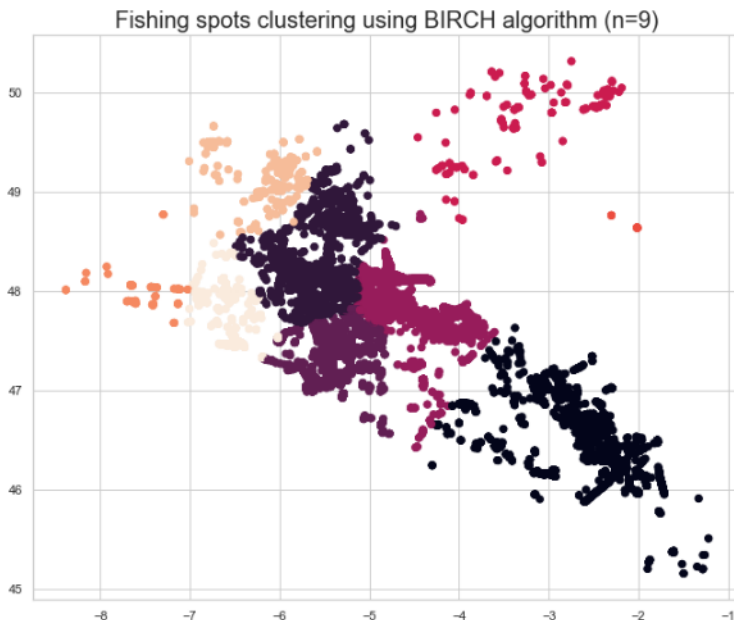


```

1 #Using BIRCH
2
3 brc = Birch(n_clusters=9)
4 clusters = brc.fit_predict(fishing_coords)
5 print('Nof clusters:', len(np.unique(clusters)))
6 clusterd_fishing_coords=fishing_coords.copy()
7 clusterd_fishing_coords['clusters']=clusters
8 plt.figure(figsize=(11,9))
9 plt.title('Fishing spots clustering using BIRCH algorithm (n=9)', fontsize=20)
10 plt.scatter(clusterd_fishing_coords.lon, clusterd_fishing_coords.lat, c=clusterd_fishing_coords['clusters'])
11 plt.show()

```

Nof clusters: 9



Παράνομη αλιεία

Ο εντοπισμός περιπτώσεων όπου υπάρχει αλιευτική δραστηριότητα σε απαγορευμένες περιοχές ανάγεται σε γεωχωρικά, χρονικώς εξαρτώμενα, queries. Εδώ γίνεται και η μεγαλύτερη αξιοποίηση των ευρετηρίων καθώς το compound ευρετήριο που δημιουργήσαμε είναι συνδυασμός χρόνου και τοποθεσίας. Τις περιοχές όπου υπάρχει απαγόρευση τις έχουμε ήδη εισάγει στη βάση σε μορφή GeoJSON με πολύγωνα. Έτσι, κάνουμε χρήση της λειτουργίας **\$geoWithin** της MongoDB, με την οποία εξετάζουμε εάν ένα σημείο βρίσκεται μέσα στα πολύγωνα που μας ενδιαφέρουν, δηλαδή στις απαγορευμένες περιοχές. Παρακάτω φαίνεται ο κώδικας και τα αντίστοιχα πλοία που εντοπίστηκαν.

Camaret area, north

```
1 import datetime
2
3 #Camaret area, north
4 start=datetime.datetime(2015,2,5,0,0,0)
5 end=datetime.datetime(2016,1,21,0,0,0)

1 camaret=db.forbidden_fishing_areas.find({"id" : 'Camaret'})

1 camaret_violations=db.nari_dynamic.find({"navigationalstatus":{"$eq":7},
2                                          "t":{"$gte":start,"$lte":end},
3                                          "location":{"$geoWithin":{"$geometry":camaret[0]}}})

1 results3=[]
2 for i in camaret_violations:
3     results3.append(i)
4
5 res3=pd.DataFrame(results3)
6 res3['sourceemsi']=res3['sourceemsi'].astype(str)

1 print('There are',len(res3['sourceemsi'].unique()),'vessels that were fishing illegally at Camaret area:')
2 print(res3['sourceemsi'].unique())

There are 2 vessels that were fishing illegally at Camaret area:
['228355800' '228853000']
```

Concarneau area, south

```
1 start=datetime.datetime(2015,9,3,0,0,0)
2 end=datetime.datetime(2015,10,15,0,0,0)
3
4 import re
5 regx = re.compile("^Concarneau", re.IGNORECASE)

1 concarneau=db.forbidden_fishing_areas.find({"id" : regx})

1 concarneau_areas=list(concarneau)

1 concarneau_violations=[]
2
3 for area in concarneau_areas:
4     violation = db.nari_dynamic.find({"navigationalstatus":{"$eq":7},
5                                     "t":{"$gte":start,"$lte":end},
6                                     "location":{"$geoWithin":{"$geometry":area}}})
7     concarneau_violations.append(violation)

1 results4=[]
2 for violations in concarneau_violations:
3     for j in violations:
4         results4.append(j)
5
6 res4=pd.DataFrame(results4)
7 res4['sourceemsi']=res4['sourceemsi'].astype(str)

1 print('There are',len(res4['sourceemsi'].unique()),'vessels that were fishing illegally at Concarneau area:')
2 print(res4['sourceemsi'].unique())

There are 2 vessels that were fishing illegally at Concarneau area:
['228827000' '227577000']
```


Συμπεράσματα

Η χρήση της μη-σχεσιακή βάσης δεδομένων MongoDB ήταν ιδιαίτερα χρήσιμη για την αποθήκευση των δεδομένων και την ανάκτησή τους. Με τη χρήση ευρετηρίων τα queries τρέχουν άμεσα και η δυνατότητα για επεξεργασία γεωχωρικών δεδομένων (GeoJSON) από τη MongoDB είναι εξαιρετικά χρήσιμη για παρόμοια προβλήματα με το δικό μας.

Για τη συσταδοποίηση τροχιών οι αλγόριθμοι που χρησιμοποιήσαμε (κυρίως density-based) προσφέρουν ικανοποιητικά αποτελέσματα και τρέχουν σε μικρό χρόνο. Επίσης η επιλογή της μετρικής του Hausdorff προσφέρει επιπλέον ταχύτητα στους υπολογισμούς.

Για τη συσταδοποίηση των σημείων, χρησιμοποιήσαμε πιο απλούς αλγορίθμους (KMeans, Birch) καθώς οι υπόλοιποι απαιτούσαν ιδιαίτερα μεγάλο χρόνο εκτέλεσης, ενώ ταυτόχρονα αναγνώριζαν πολύ μεγάλο πλήθος συστάδων.

Για την εξέταση περιπτώσεων παράνομης αλιευτικής δραστηριότητας, το γεγονός ότι έχουμε compound ευρετήριο ως προς χρόνο και τόπο ταυτόχρονα, μας προσφέρει ταχύτατα αποτελέσματα. Εντοπίσαμε παράνομη δραστηριότητα και στις δύο περιοχές όπου υπήρξε απαγόρευση, για διαφορετικά χρονικά διαστήματα.