

## AIS Anomaly Detection

- Fetch Broadcast Data per Year
- Iterate thru each unique MMSI
  - Identify anomal coordinate for a given vessel
- - Flag any Anomaly and Aggregate all Vessel Data into dataframe along with various stats: Mean, Median, etc
- Output to Statistic\_{Year}.csv

```
In [1]: # from IPython.display import Image, HTML
import os
import numpy as np
import math
import pandas as pd
import datetime
from glob import glob
import geopy.distance
import folium
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
import seaborn as sns; sns.set()

import warnings
warnings.filterwarnings("ignore")           # Suppress Warning
```

## Global Variables

```
In [2]: # s3://vault-data-corpus/vessel data/ConsolidatedAIS/
WorkingFolder = "data/vessel data/ConsolidatedAIS/"
OutputDir = WorkingFolder

PROC_YEAR = '2009'
MAX_CLUSTER = 5
MIN_PROCESS_ROW = 500           # Require min number of rows available to calc cluster
```

## Load Broadcast Data

```
In [3]: Broadcast = pd.read_csv(WorkingFolder + "Broadcast_{}.csv".format(PROC_YEAR), sep=",", parse_dates=['date_time'])
Broadcast.head()
```

```
Out[3]:
```

	mmsi_id	date_time	lat	lon	speed_over_ground	course_over_ground	voyage_id	heading	status
0	367047170	2008-12-31 23:58:59	45.633627	-122.715003	0	199	1	511	0
1	366763770	2008-12-31 23:58:59	46.027173	-122.869968	0	325	2	511	0
2	368494000	2008-12-31 23:58:59	47.542502	-122.330622	0	354	3	511	0
3	366116000	2008-12-31 23:58:59	48.512612	-122.645723	0	345	4	55	0
4	316003289	2008-12-31 23:58:59	49.144225	-123.033517	0	172	5	96	15

```
In [4]: print("Raw Count:", Broadcast.shape[0])
```

Raw Count: 17749932

```
In [5]: # Make sure voyage id is not null
Broadcast['voyage_id'] = Broadcast['voyage_id'].fillna(0)
Broadcast = Broadcast.astype({"voyage_id": int})           # cast type to int
```

## Clustering via K-means

- <https://github.com/JosephMagiya/Clustering-GPS-Co-ordinates--Forming-Regions/blob/master/Clustering-GPS-Co-ordinates--Forming-Regions.ipynb> (<https://github.com/JosephMagiya/Clustering-GPS-Co-ordinates--Forming-Regions/blob/master/Clustering-GPS-Co-ordinates--Forming-Regions.ipynb>)



```
In [8]: Stat
```

Out[8]:

	date_time	lat	lon	speed_over_ground	course_over_ground	PingDate
count	791	791.000000	791.000000	791.000000	791.000000	791
max	2009-01-31 23:58:59	32.798445	-122.471355	19.000000	127.000000	2009-01-31
min	2009-01-31 10:34:00	30.910263	-125.996375	8.000000	118.000000	2009-01-31
mean	2009-01-31 17:16:06.317319936	31.835202	-124.190095	15.906448	122.116308	NaN
median	NaN	31.761025	-124.049887	18.000000	122.000000	NaN
std	NaN	0.519859	0.970443	3.685256	1.470081	NaN

```
In [9]: df.head()
```

Out[9]:

	date_time	lat	lon	speed_over_ground	course_over_ground
17749646	2009-01-31 23:58:59	34.242787	-120.000647	15	284