



HACKtheMACHINE

September 21-23rd

Team: #WestCoastIsTheBestCoast

Joe Ortiz, James Eitelberg, Chris Chen, Jessica Kimball, Chris Cirullo, Thuy Nguyen-Cirullo (SSC PAC); Richard Brereton (COMNAVSURFPAC);
Rock Pereira (Consultant)

Track 2: Data Science and The Seven Seas

#WestCoastIsBestCoast Approach #1

- Clean-Up and Prioritization
 - Remove unneeded data (8 fields, e.g. IMO, width)
 - Filter erroneous data (e.g. Tug Boats, negative COG, moored, anchored vessels)
 - Divide data into 10-minute periods
- Processing and Visualization
 - Determine distance between vessels using k-d trees
 - Import data into ArcGIS map over maritime layers

Results

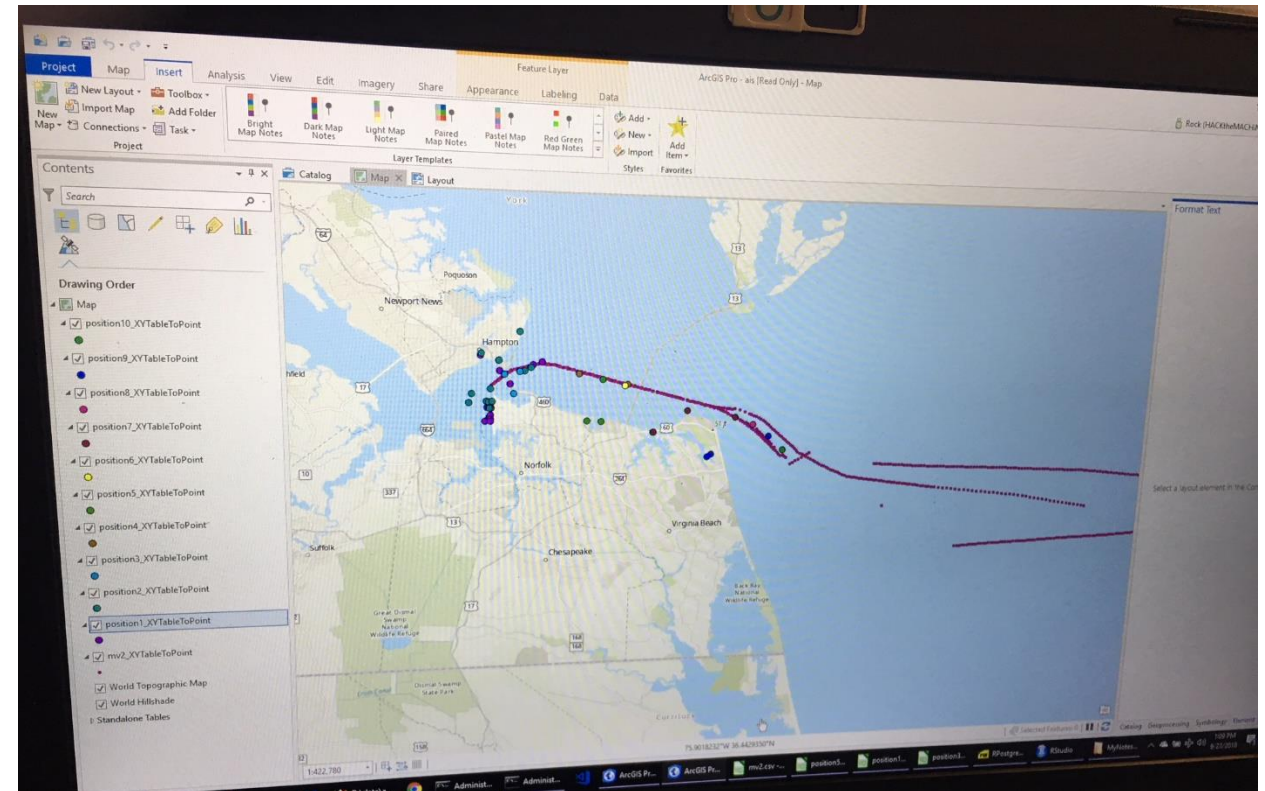
- Computation time based on same data set:
 - Brute force: ~30/40 min
- Our machine learning based kd-tree algorithm we were able to improve performance to 10-15 min
 - > 50% improvement

#WestCoastIsBestCoast Approach #2

Simple geospatial calculation filtered by ship types
Follow that ship and create view frames

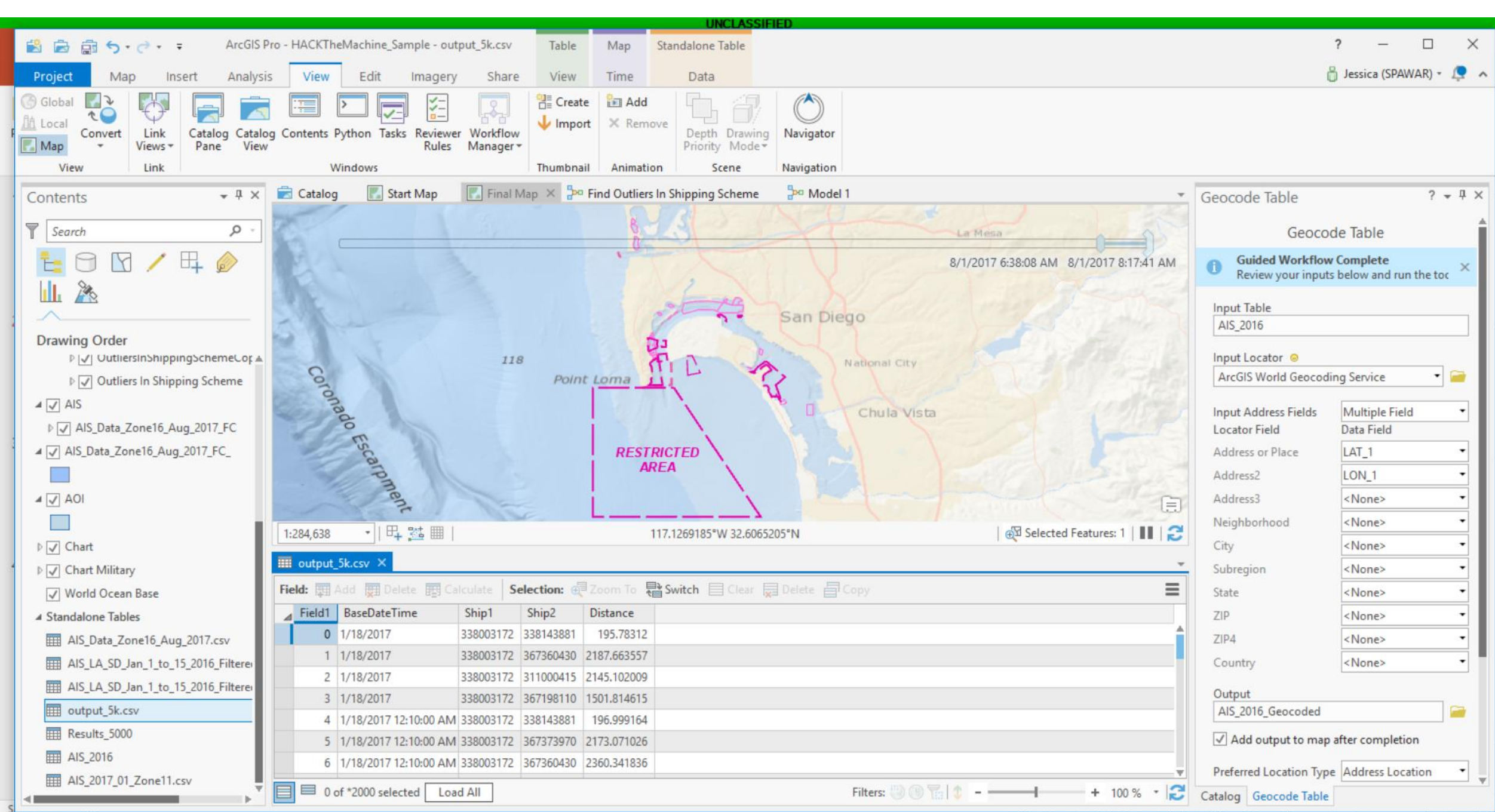
AIS data stored into AWS Cloud!

- Data for 3 years in Zone 10
- All zones for 1 month (DEC 2017)

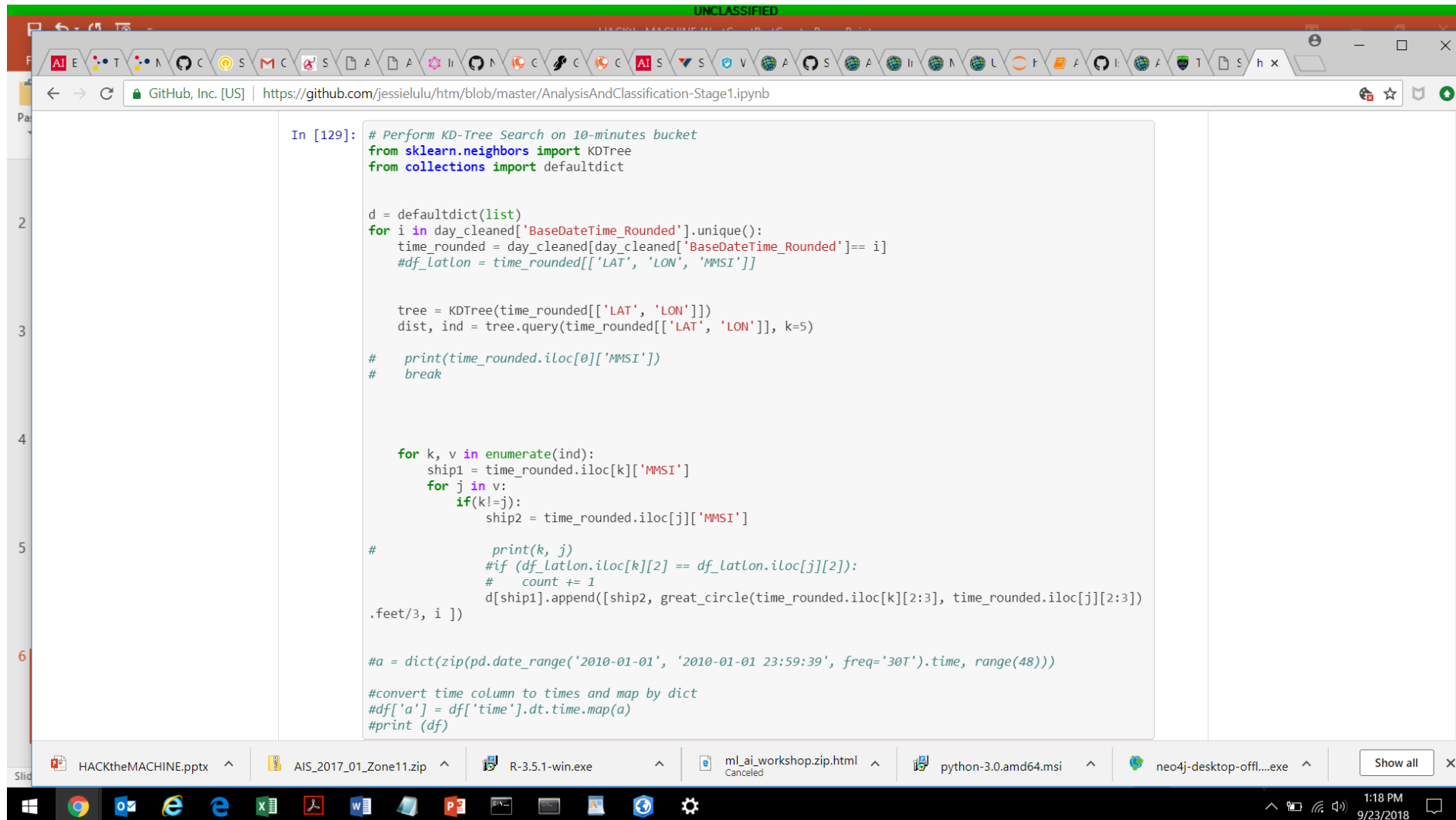


Our Approach Challenge #2 - Continue

- Stage 2:
 - Preprocess each of the features:
 - For continuous variables, normalize it by subtract it by the means and divided by standard deviation
 - For categorical variables, encode it to a number
 - PCA and/or Factorial Analysis
 - To determine which features are the most important ones.
 - Multiclass Classification Algorithms.
 - ie: K-means Clustering, and etc.



Jupyter Notebook Example



```
In [129]: # Perform KD-Tree Search on 10-minutes bucket
from sklearn.neighbors import KDTree
from collections import defaultdict

d = defaultdict(list)
for i in day_cleaned['BaseDateTime_Rounded'].unique():
    time_rounded = day_cleaned[day_cleaned['BaseDateTime_Rounded']== i]
    #df_latlon = time_rounded[['LAT', 'LON', 'MMSI']]

    tree = KDTree(time_rounded[['LAT', 'LON']])
    dist, ind = tree.query(time_rounded[['LAT', 'LON']], k=5)

    # print(time_rounded.iloc[0]['MMSI'])
    # break

    for k, v in enumerate(ind):
        ship1 = time_rounded.iloc[k]['MMSI']
        for j in v:
            if(k!=j):
                ship2 = time_rounded.iloc[j]['MMSI']

            # print(k, j)
            #if (df_latlon.iloc[k][2] == df_latlon.iloc[j][2]):
            #    count += 1
            d[ship1].append([ship2, great_circle(time_rounded.iloc[k][2:3], time_rounded.iloc[j][2:3])
            .feet/3, i ])

#a = dict(zip(pd.date_range('2010-01-01', '2010-01-01 23:59:39', freq='30T').time, range(48)))

#convert time column to times and map by dict
#df['a'] = df['time'].dt.time.map(a)
#print (df)
```