

All code, plots reside here:

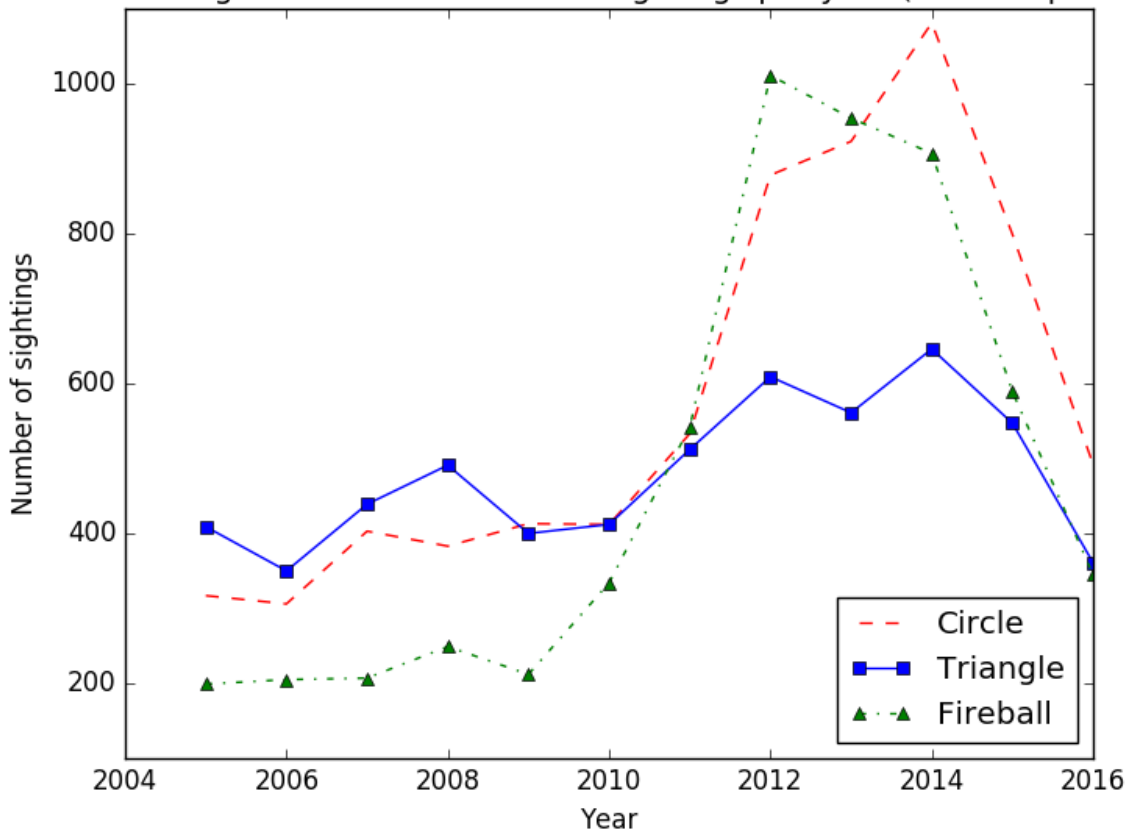
<https://github.com/karunmj/usu-coursework/tree/master/cs5660datasc/hw/hw2>

Data collection involved using request and beautiful soup libraries to get UFO data for circles, triangles and fireball shape sightings in the form of a Pandas data frame object in Python. The column labels include date of sighting, time of sighting, city, state, shape, duration, summary and posted date.

Data preprocessing involved including sighting between 1/1/2005 and 9/22/2016, representing duration of sighting in seconds and sightings in either of the 50 (+1 by including DC) states of US.

1. Box plot of duration of UFO sightings of each shape  
(not done)
2. Time series figure with the number of sightings per year (one line per shape)

Time series figure with the number of sightings per year (one line per shape)



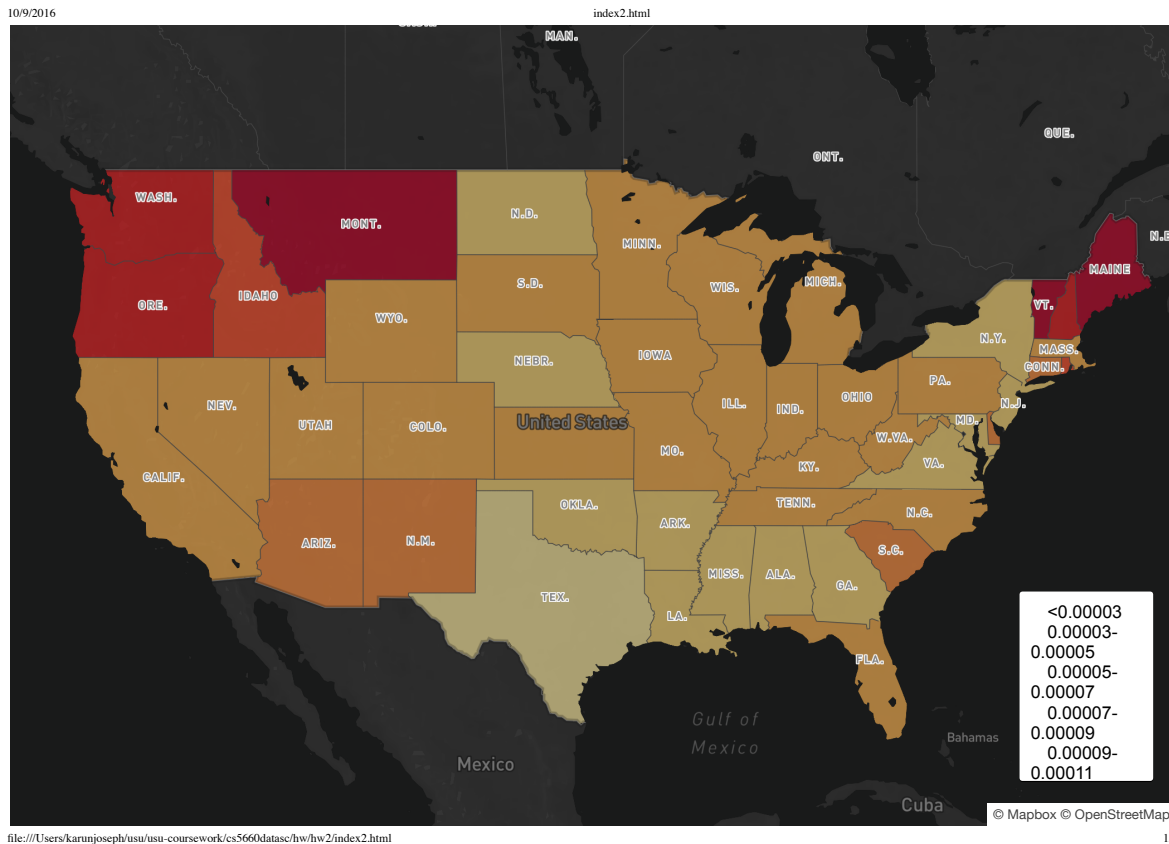
There seems to be higher number of sightings during 2012 to 2014 than the rest of years.



CS 5665 Homework #2  
Karun Joseph, A02240287

State population numbers were obtained from US Census API service for

4. Normalized sightings by state population visualized on a map  
Mapbox Studio had been used to generate a static HTML page with the normalized sightings.



It looks like states like Montana, Vermont and Maine had the highest number of sightings. These states also have very low population, hence higher sightings per person. Also, these states have fewer cities that result in lower light pollution levels. Residents are more likely to see moving planetary objects like comets, meteors which could be mistaken for UFOs.

For classification, Python's sklearn tree package had been used to build the decision tree classifier. Since this package doesn't deal with categorical type of attribute directly [which is in our case], they had to be vectorized. For example, given an observation's attributes in training data set

| time of day | region |
|-------------|--------|
| night       | s      |

it is converted to:

| region<br>=mw | region<br>=ne | region<br>=s | region<br>=w | time_of_day=<br>afternoon | time_of_day=<br>evening | time_of_day=<br>morning | time_of_da<br>y=night |
|---------------|---------------|--------------|--------------|---------------------------|-------------------------|-------------------------|-----------------------|
| 0             | 0             | 1            | 0            | 0                         | 0                       | 0                       | 1                     |

CS 5665 Homework #2  
Karun Joseph, A02240287

Similarly, the class labels ‘Circle’, ‘Fireball’ and ‘Triangle’ were mapped to 0, 1 and 2 respectively. Gini criterion had been used to build the classifier.

5. Classification accuracy of decision tree using test set data  
Python’s sklearn metrics’ precision score had been used to calculate the ratio  $\frac{t_p}{t_p + f_p}$ , where  $t_p$  is the number of true positives and  $f_p$  the number of false positives, for each class

| Circle     | Fireball   | Triangle   |
|------------|------------|------------|
| 0.43863816 | 0.35789094 | 0.33474576 |

The accuracy from sklearn metrics’ accuracy\_score is 0.391364661002

The confusion matrix of predicted vs true class labels of test set data were as follows

| True label | Circle   | 1108            | 1067     | 197      |
|------------|----------|-----------------|----------|----------|
|            | Fireball | 733             | 991      | 117      |
|            | Triangle | 685             | 711      | 158      |
|            |          | Circle          | Fireball | Triangle |
|            |          | Predicted label |          |          |

Overall, the decision tree classifier had done a poor job.

6. Illustration of built decision tree

