# Prediction of West Nile virus in Chicago

## LAL BABU SAH

## MAY 30, 2020

### Introduction:

In this project separate train, test, weather, and spray dataset are given. And I have to predict whether or not 'West Nile virus' in Chicago is present, for a given time, location, and species using given dataset. Here, I will use the Classification machine learning algorithm for the prediction of 'West Nile virus'.

### Data Acquisition:

Read the train, test, weather, and spray dataset using pandas.

### Data Wrangling:

First merged 'train' dataset with 'weather' dataset along the common column (Primary Key) 'Date'. Then I have matched the 'Date' and 'location' of the GIS dataset with the merged dataset and added the new column(spray) in the merged dataset (added spray is equal to 1 if matched else 0).

After that, I have dropped some columns which are not important for prediction, like from ('Block', 'Address', 'Address Number and Street') these all contains address, dropped 'Address', 'Address Number and Street' but not 'Block'.

Column ('species') contain different types of species name which are the object type, I have converted these in numerical data. Some columns contain the numerical data, but it is the object type, converted these into 'int', or 'float' type. Some columns ('StnPressure', 'SnowFall') contain 'M' which are missing data, replaced it with '0'.

### Data Exploration:

Calculated the Pearson Correlation Coefficient and the p-value of different columns and dropped the columns which are not important for prediction on the basis of correlation and p-value. Checked the outliers of different columns, to remove the outliers if present.

### Model Development:

After data exploration normalized the dataset and then split into the train and test dataset. Here I have used classification machine learning algorithm because target ('WnvPresent') is categorical ('Yes', 'No'). After that, preprocess the test dataset and predicted the 'WnvPresent'.

**Data Visualization using Tableau:**

In tableau uploaded the merged train dataset for the dashboard preparation.
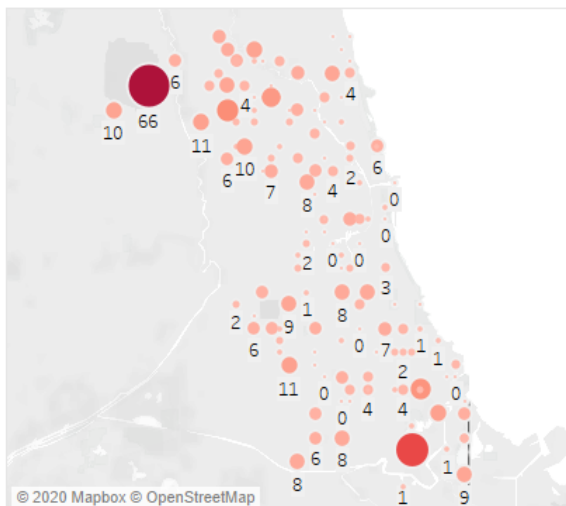
Using geographical coordinates plotted map, which is showing in which location 'West Nile virus' is present more or less for different years.

Using the columns 'Date' and 'WnvPresent' plotted year wise trend, which shows in which month and year 'West Nile virus' is present more or less.
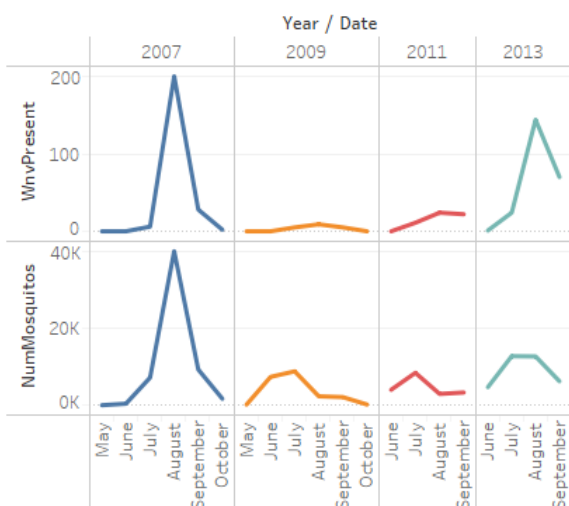
Using columns 'Block', 'Street', 'NumMosquitos', and 'WnvPresent' plotted bar graph, which shows the number of 'NumMosquitos' and 'WnvPresent' for a given 'Block' and 'Street' for different years.

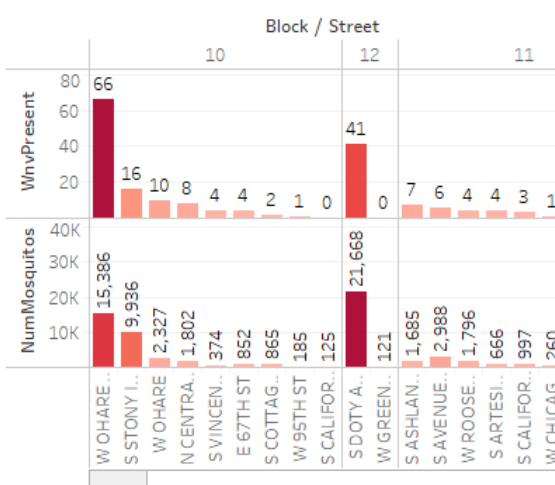Using columns 'Species' and 'WnvPresent' plotted graph, which shows the number of 'WnvPresent' for different species for different years.