





# West Nile Virus Prediction

Audrey, Akhila, Shaun



# Table of Contents

- 1) Problem Statement
  - 2) Data Cleaning and EDA
  - 3) Preprocessing
  - 4) Modelling
  - 5) Limitations & Recommendations
  - 6) Cost-Benefit Analysis
  - 7) Conclusion & Recommendations
-

# The West Nile Virus

- Leading cause of mosquito-borne disease
- Potentially fatal
- No Vaccine



# Problem Statement

# Problem Statement

## Project Aim

- Predicting *West Nile Virus* (WNV) in Chicago
- Utilising weather features

## Purpose

- Assist CDC and CPHD in combating WNV

# Data Cleaning

# Train Test

- 1) **Incorrect dtype**
- 2) **Drop (Redundant)**
- 3) **Drop (Duplicates)**
- 4) **Drop (Not in Test)**



# Spray

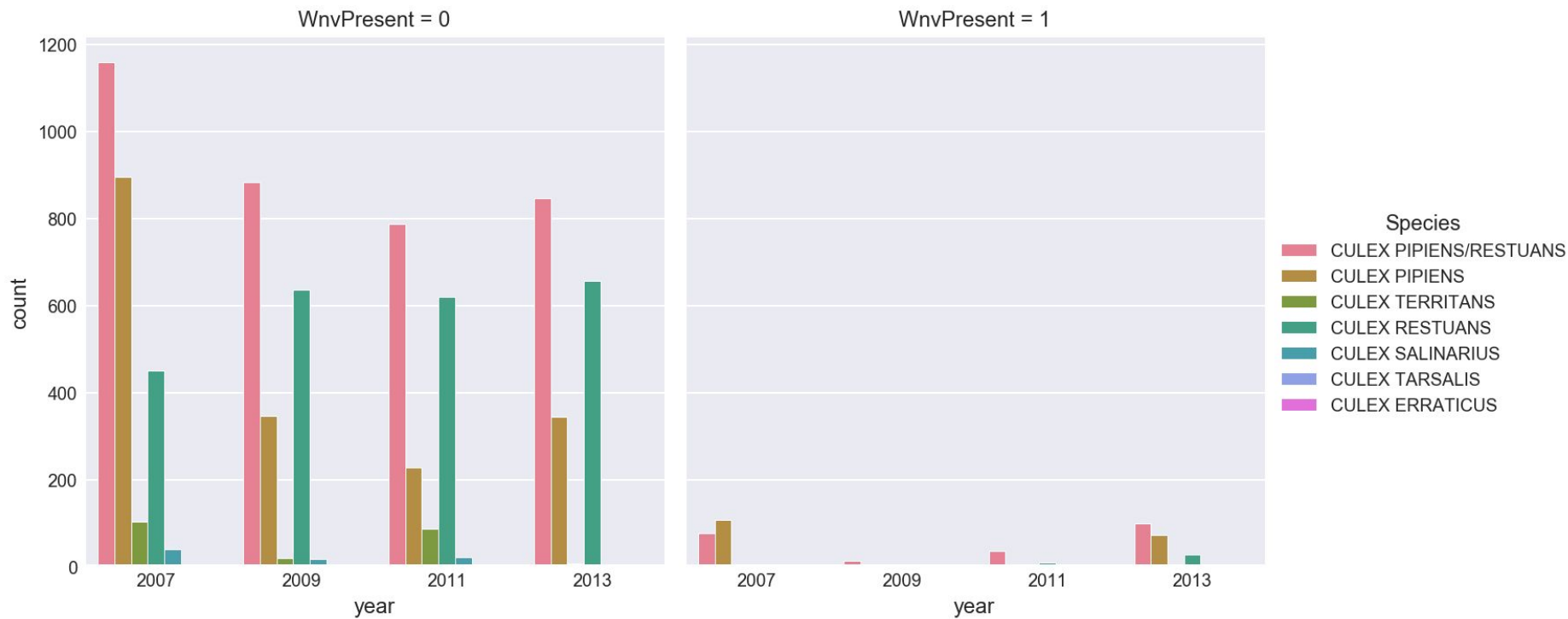
- 1) **Incorrect dtype**
- 2) **Drop (Redundant)**
- 3) **Drop (Duplicates)**

# Weather

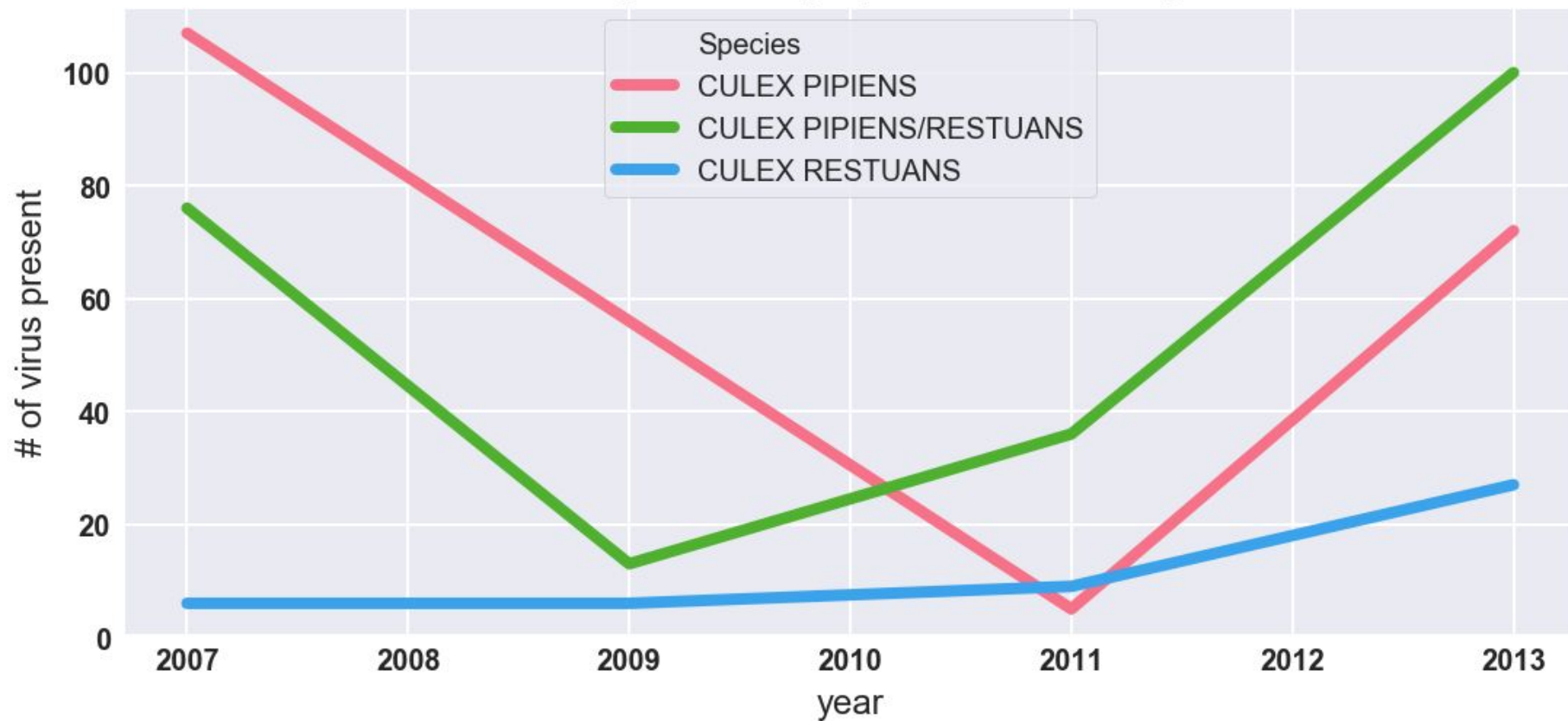
- 1) Incorrect dtype
- 2) Drop (Station 2)
- 3) Drop (Missing Info)
- 4) Impute (“M”, “T”)



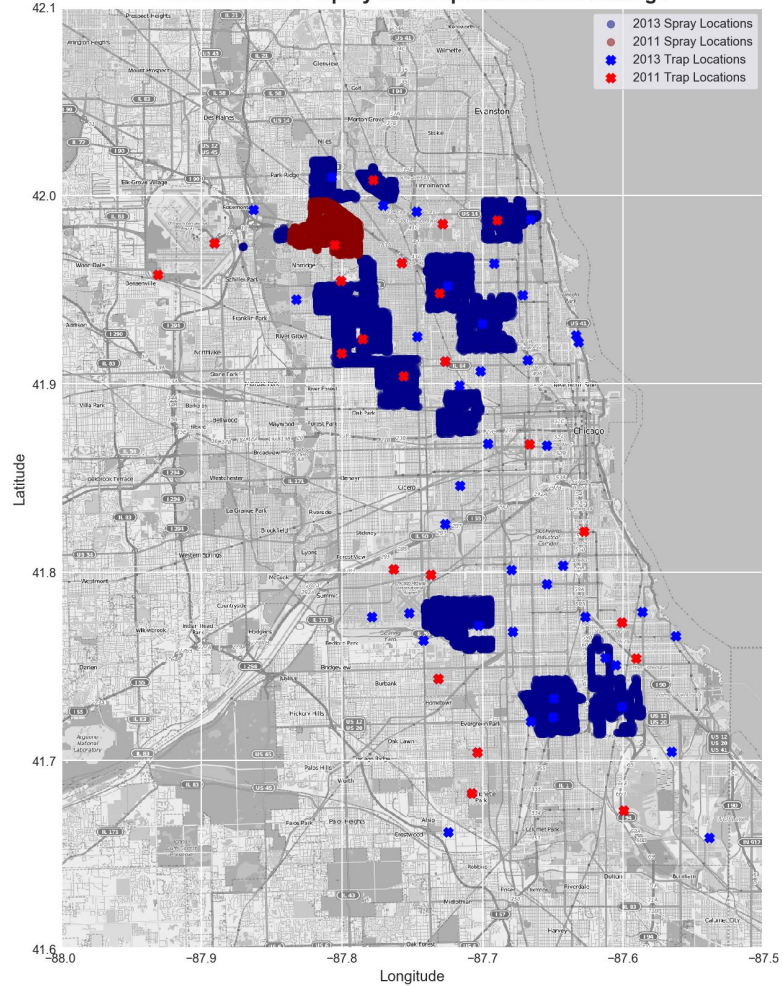
# EDA



**# of virus present by species over the years**



Map of Chicago showing 2013 Spray Locations (blue dots), 2011 Spray Locations (red dots), 2013 Trap Locations (blue stars), and 2011 Trap Locations (red stars). The map includes latitude and longitude coordinates and a legend.



# Preprocessing

# Preprocessing

## Features Engineered

1. Dummy variables for species created in both train and test set
2. Parsed dates into year, month, week of year, month of year
3. Coded hot and wet conditions based on dew point and average temperature
4. Merged weather dataset to train set for modeling on Date

## Dropped

1. Spray set is omitted for now due to the lack of information over the years.
2. Temperature related & Dewpoint features



# Modelling

# Models Tried

1. Logistic Regression
2. Gradient Boost Classifier
3. Gradient Boost Classifier with GridSearchcv
4. Random Forest Classifier
5. XGBoost Classifier
6. Decision Trees

# Modeling Results

Model	Hyper-parameters	Train Set AUC Score	Test Set AUC Score	Kaggle Score
Logistic Regression	-	<b>0.74</b>	<b>0.78</b>	<b>0.50</b>
Gradient Boost Classifier	-	<b>0.90</b>	<b>0.87</b>	<b>0.56</b>
	'learning_rate': 0.08, 'max_depth': 2, 'n_estimators': 100	<b>0.87</b>	<b>0.87</b>	<b>0.49</b>
Random Forest Classifier	-	<b>0.99</b>	<b>0.78</b>	<b>0.52</b>
XGBoost Classifier	n_estimators = 500	<b>0.96</b>	<b>0.87</b>	<b>0.53</b>
Decision Trees	-	<b>0.99</b>	<b>0.58</b>	<b>0.52</b>

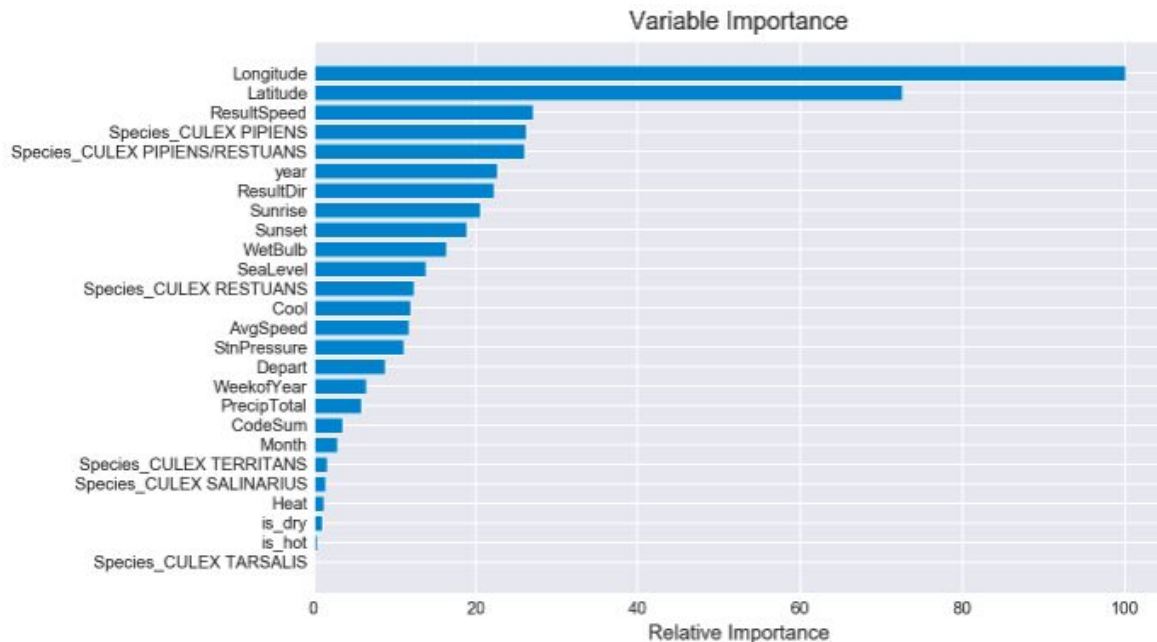
# Model Evaluation

Baseline score vs. Chosen model score

Model	Train Set	Test Set
Logistic Regression	0.74	0.78
Gradient Boost Classifier	0.90	0.87

Heavily imbalanced dataset

Important Features



# Insights & Conclusions

- 1) Least variance between test and training scores >> not overfit
- 2) Feature importance suggests the prevalence of WN virus in certain locations for specific species
- 3) Time of the year with longer days influences virus presence
- 4) Features highlighted to be used to determine where and when to focus for spraying efforts



# Limitations & Recommendations



# Limitations & Recommendations

## LIMITATION

Model scored significantly lower on out of sample Kaggle set  
>> *probably due method of creating internal validation set vs. Kaggle's approach*

## RECOMMENDATIONS

Elaborate feature engineering with time series functions on weather dataset

Include dummified variable for top traps

*NumMosquito* feature even with its high correlation to virus presence had to ignored

Availability of spray information for more years would help the model score better



# Cost-Benefit Analysis





# Cost-Benefit Analysis

	<b>Spraying</b>	<b>Not Spraying</b>
<b>Costs</b>	<ul style="list-style-type: none"><li>● Chemical spray (Zenivex)</li><li>● Labour work<ul style="list-style-type: none"><li>○ Approx. \$0.75 per acre of spray</li></ul></li></ul>	<ul style="list-style-type: none"><li>● Medical cost<ul style="list-style-type: none"><li>○ \$33,143 per inpatient</li><li>○ \$6,317 per outpatient</li><li>○ \$18,097 per patient spent time in a nursing home</li></ul></li><li>● Productivity loss<ul style="list-style-type: none"><li>○ \$58,935 per personal income</li></ul></li></ul>

Source:

<http://www.gfmosquito.com/wp-content/uploads/2013/06/2013-North-Dakota-Bid-Tabulation.pdf>

<https://www.statista.com/statistics/205235/per-capita-personal-income-in-illinois/>

<https://www.cdc.gov/westnile/resources/pdfs/data/WNV-Disease-Cases-PVDs-by-State-2018-P.pdf>

# Cost-Benefit Analysis

	COST	BENEFIT
Spraying	<ul style="list-style-type: none"><li>Vector control cost</li></ul> <p>Total = ~ \$144K per person</p>	Human Life Saved
Not Spraying	Human Life Loss	<ul style="list-style-type: none"><li>Medical bills (~\$46,530 per person)</li><li>Productivity loss ~\$58,935 per person)</li></ul> <p>Total = ~\$108K per person</p>

**Total \$ Cost > Benefits ?**

# Cost-Benefit Analysis

	<b>Spraying</b>	<b>Not Spraying</b>
<b>Benefits</b>	<ul style="list-style-type: none"><li>● <b>Human Life Saved</b><ul style="list-style-type: none"><li>○ Improved quality of life</li><li>○ Increased workplace productivity</li><li>○ Savings in hospital bills</li><li>○ Attract visitors → economic benefits</li></ul></li></ul>	<ul style="list-style-type: none"><li>● <b>Human Life Loss</b><ul style="list-style-type: none"><li>○ Long term/Wider mental health issues</li><li>○ Lower quality of workforce and reduce output productivity</li><li>○ Negative impact on tourism and entertainment sectors</li></ul></li></ul>

**Human Benefits of spraying outweighs the Costs!**

# Cost-Benefit Analysis

<b>Year</b>	<b>Acre coverage*</b>	<b>Projected Annual Costs @ \$0.75 per Acre</b>
<b>2020</b>	~16m	~\$12m
<b>2021</b>	~14.4m	~\$10.8m
<b>2022</b>	~13m	~\$9.75m

\*Source:

<https://www.cdc.gov/westnile/statsmaps/cumMapsData.html>

<https://www.cdc.gov/westnile/vectorcontrol/aerial-spraying.html>

(assumption: 10% decrease in levels of pesticide coverage annually)



# Conclusion and Recommendations



# Conclusion & Recommendations

Since humans are gregarious and following CDC's mission statement, we must think of a regional approach to WNV vector control..... Spray!

- Consider strategically spraying at locations with highest infections
- Conduct spraying during hotter months i.e. August and September
- Educate and promote the public to:
  - Use insect repellent
  - Wear long-sleeved shirts and pants
  - Take steps to control mosquitoes indoors and outdoors
    - I.e. Remove standing water where mosquitoes could lay eggs



Questions?

