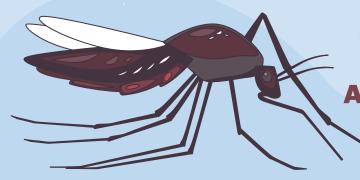# Project 4: West Nile Virus Prediction

DSI-28
10 June 2022
Adi, Calvin, Joel, Priscilla, Yong Lim

# Agenda

# 1

## Introduction and Problem Statement

# Problem Statement



To prevent a WNV outbreak in Chicago, the Chicago Department of Public Health (CDPH) has tasked its data science team to develop a _predictive model to detect areas highly likely to have WNV_.

In addition, CDPH has requested for our expertise in advising the _areas and timings to spray pesticide_ as part of the terms in the contract. The advice also includes data-driven _benefits_ of spraying, and annual _costs_ estimates to be used in their price negotiations for the new spraying contract.

# 2

# Data Cleaning

# Data Provided

| Data | Description | Timeframe |
|------|-------------|-----------|
| Train | Data set to train predictive model | 29 May 2007 to 26 Sep 2013 |
| Test | Data set to produce prediction results | 11 June 2008 to 2 Oct 2014 |
| Spray | Time, date and location of previous sprays | 29 Aug 2011 to 5 Sep 2013 |
| Weather | Meteorological information of Chicago | 1 May 2007 to 31 Oct 2014 |

# Data Cleaning - Train Dataset

## Removing duplicate records

Based on Date, Species and Trap

# 3

# Exploratory Data Analysis (EDA)

# Mosquito numbers trend closely with temperature over the year



Mosquito Count and Average Temperature (Monthly)

- Mosquito count increases with temperature, peaking in July-August

- Mosquito count decreases in later months as temperature drops

- Slight lag between mosquito count and temperature

# Mosquito count tend to peak after rainfall



Peak in rainfall followed shortly by peak in mosquito numbers

# Mosquito count increase with daytime length



Effect of Day Time on Mosquitos Trapped

More mosquitoes are detected as day length increases

# Spraying in response to WNV cases, as opposed to mosquito numbers



- Spraying counts increased almost instantaneously to emergence of WNV cases in July. (left)
- Increase in number of mosquitoes started much earlier in May. (right)
- Spraying may have started only in response to receiving reports of WNV cases

# Summary of Observations

Mosquito count and Virus count are increased with:

- Higher temperatures
- Higher rainfall/ precipitation
- Longer daytime

And are effectively lowered by:

- Spraying insecticide

# Culex Pipiens and Restuans species carry WNV



West Nile Virus found exclusively in:
- Culex Pipiens
- Culex Restuans

# Disparity in spraying and actual WNV outbreak locations



July

August

September

- Areas sprayed with insecticide (green) and locations with WNV cases (red)

- Sprayed area tend to cover much larger area than the neighbourhood - wastage

- Problem areas not receiving spraying

# 4

## Feature Engineering

# Factors Affecting Mosquito Breeding

## Environment

Temperature, Relative Humidity, Precipitation / Rain, Wind

## Time

Week, Night / Day, 7 days lag

## Location

Latitude, Longitude , WNV risk

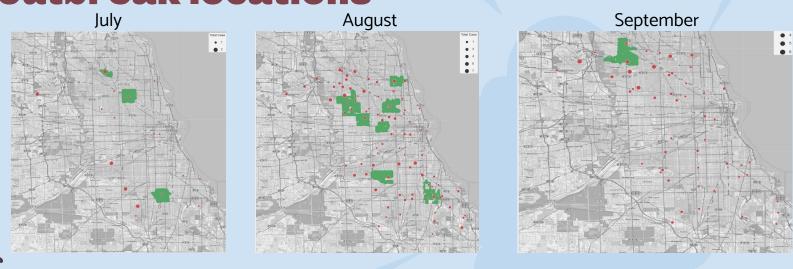## Species

Culex Pipens, Culex Restuan

*Source 1*: Predicting Culex pipiens/restuans population dynamics by interval lagged weather data
*Source 2*: When are mosquitos most active

# Environment Added Features



Relative Humidity

$$RH = 100 \times \left[ \dfrac{e^{\frac{17.625 \times D_p}{243.04 + D_p}}}{e^{\frac{17.625 \times T}{243.04 + T}}} \right]$$

Dp – Dewpoint Temp
T – Ave Temp

# Environment Dataset

Average Temperature
Dew Point
Precipitation
Wind Speed / Direction
Station Pressure
Sea Level
Rain / Thunderstorm / Mist

**Location Added Features**

WNV Risk
Low – 0 to 2 WNV cases
Medium – 3 to 5 WNV cases
High – Above 5 WNV cases

**Location Dataset**

Latitude
Longitude

# Time
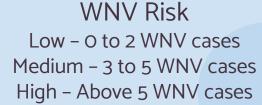## Added Features

# Species

Night Time
Day Time
Week
7 Days Lag for temp,
dewpoint & precipitation

Culex Pipens
Culex Restuan

# 5
# Modelling

Logistic Regression
Random Forest
AdaBoost
Gradient Boost
XGBoost
Support Vector Machine
K-Nearest Neighbours

# Model Evaluation Metric - AUC ROC



- Evaluates how good the model distinguish the classes

- Higher AUC -> Better in predicting 0 class and 1 class

- Good Model -> Accurately predict the presence of the virus

Source of illustration:
https://towardsdatascience.com/understanding-auc-roc-curve-68b2303cc9c5

# Logistic Regression



Used as baseline model

**ROC AUC (Train)**

0.7771

**ROC AUC (Validation)**

0.8554

# Random Forest

Training Dataset

$\Sigma$

Final Result

## ROC AUC (Train)
0.9169

## ROC AUC (Validation)
0.8790

# Adaptive Boosting (AdaBoost)



**ROC AUC (Train)**

0.8752

**ROC AUC (Validation)**

0.8809

# Gradient Boosting



ROC AUC
(Train)

0.9065

ROC AUC
(Validation)

0.8808

# eXtreme Gradient Boosting (XGBoost)



**ROC AUC (Train)**

0.8762

**ROC AUC (Validation)**

0.8799

# k-Nearest Neighbours (KNN)



k= 3

k= 7

k= 11

## ROC AUC (Train)

0.9361

## ROC AUC (Validation)

0.8592

**6**

# Modelling Results

# Models Performance Summary

| Model | ROC AUC | | |
|---|---|---|---|
| | Train (A) | Validation (B) | Difference (B) - (A) |
| Logistic Regression | 0.8388 | 0.8542 | 0.0154 |
| Random Forest | 0.9169 | 0.8790 | -0.0379 |
| AdaBoost | 0.8752 | 0.8809 | 0.0057 |
| Gradient Boost | 0.9065 | 0.8808 | -0.0257 |
| Extreme Gradient Boost | 0.8762 | 0.8799 | 0.0037 |
| Support Vector Machine | 0.893 | 0.8772 | -0.0158 |
| k-Nearest Neighbours | 0.9361 | 0.8592 | -0.0769 |

# Top 10 AdaBoost Predictors



## Top 10 Predictors

| Feature | |
|---------|---|
| Week | |
| WnvRisk_medium | |
| Species_CULEX PIPIENS/RESTUANS | |
| Species_CULEX PIPIENS | |
| ResultDir | |
| StnPressure | |
| Species_CULEX RESTUANS | |
| AvgSpeed | |
| WetBulb | |
| Latitude | |

Coefficient Values

**7**

**Cost-Benefit Analysis and Recommendations**

# Cost of Mosquito Abatement Program 2023

## ~USD 520,698*

**which includes:**

- Weekly Environmental Surveillance (~ 147 gravid traps)
- Conduct Larviciding (~ 190 acres)
- Conduct Adulticiding (~ 100 miles)

* based on Contract (PO) Number 17068 **"SLE Vector Mosquito Abatement Program"** awarded to *Vector Disease Control International (VDCI)*

# Average Total Economic Cost for 2023: ~USD 2,800,100*

## Average Cost Per Person: ~USD 176,071*

Assumption based on:

- Average of people infected of 17 throughout 2012-2021 [+]

- Average death rate of 2 throughout 2014-2016 [+]

\* Forecasted from data source: "**Initial and Long-Term Costs of Patients Hospitalized with West Nile Virus Disease**" (_Source_) paper by _Centers for Disease Control and Prevention (CDC)_ dated 05 Mar 2014

[+] **Source:** West Nile Virus Surveillance Reports

# Benefits > Costs on Average by 4 times
## → Continue Mosquito Abatement Program

| | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 |
|---|---|---|---|---|---|---|---|---|---|---|
| **Human Cases** | 22 | 1 | 6 | 16 | **49** | 6 | 42 | 2 | 11 | 10 |
| **Human Cases (Fatality)** | - | - | - | 3 | 2 | 1 | - | - | - | - |

**Source:** <u>West Nile Virus Surveillance Reports</u>

**Worst case scenario** based on 2016 records, with **49** cases reported and **2** casualties,

the total lost instead would be: **~USD 8,434,382** , a whooping **16 times** from the 2023 abatement cost

# Recommendations

### Increase Larviciding Initiation

At the location which has a high risk of West Nile Virus emergence (WnV case > 5), additional 800 acres => USD 77,920
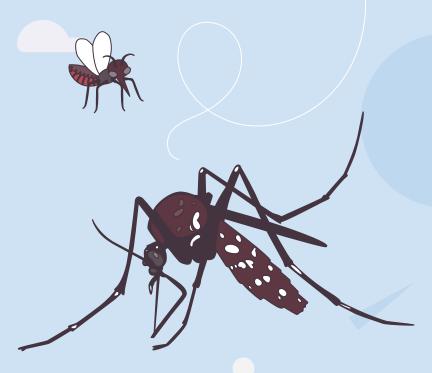
### Lower the Threshold to Activate Adulticiding

For the month of June and July, to suppress the population of mosquitoes , additional 100 miles => USD 12,060

### Conduct Awareness Roadshow

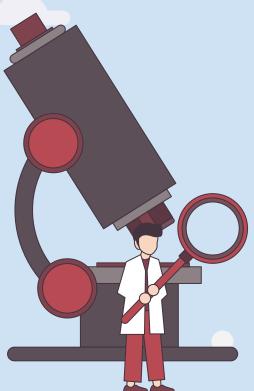Before the breeding season start. Helps to reduce potential breeding location, estimated cost ~USD 15,000

**8**

**Limitations and Future Steps**

# Limitations

### Train & Weather Dataset Range
→ Only includes data from 2007 to 2014. Weather conditions may changed since 2014
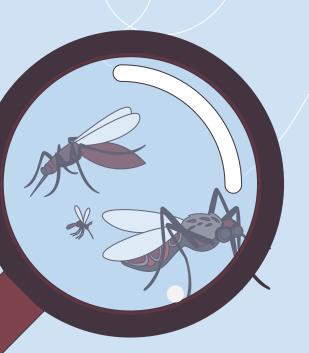
### Train Dataset Size
→ Train dataset size is comparatively smaller compare to test dataset

### Time Constraint
→ Limited time for hyperparameter tuning to obtain better performing model

# Future Works

## Effect of the New Spraying Schedule
→ To update the model with latest data and check if it is effective or the trend still persists
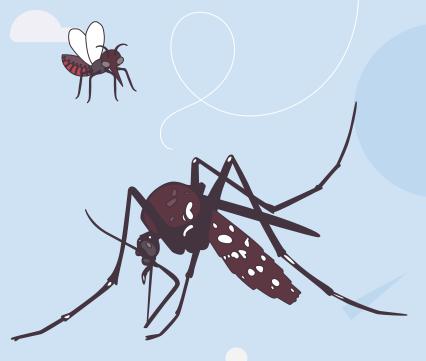
## Data on Location and No of Larvaes
→ To study the trend on the larvaes found to have a better plan on early prevention

## New Technique on Treating Features
→ PCA can be tested for treating the collinearity between existing features

# 9

# Conclusion

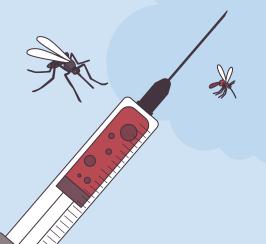# Conclusions

## Recommended Model

AdaBoost

## Estimated Costs

~USD 625,678

## Targeted Mosquito Abatement Efforts

1. Increase Larviciding Initiation at high risk areas
2. Lower the Threshold to Activate Adulticiding during June and July
3. Conduct Awareness Roadshow before May

# Thank You

**Be Ready,**
**Stay Vigilant,**
**&**
**Abolish !**