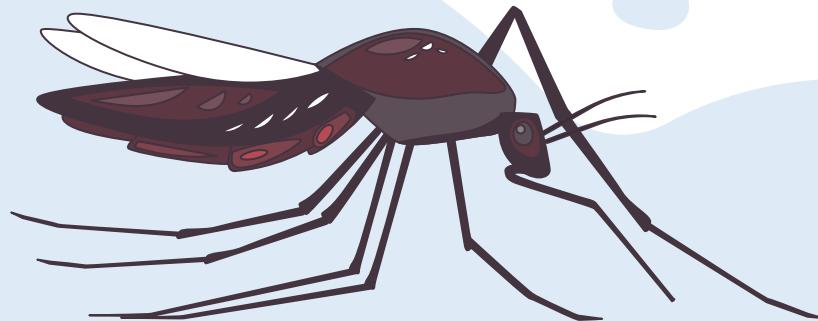




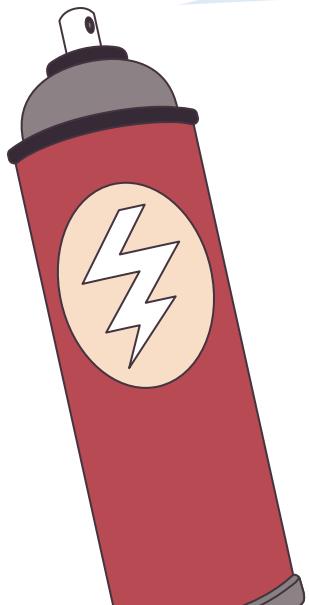
# West Nile Virus Prediction

*Big problems come in small packages*



**Members**  
Mubin  
Weihan  
Zhenming

# Contents



**01**

## **What's the big deal?**

- Background & Motivations

**02**

## **What do we know?**

- Data Collection & Exploration
- WNV Modelling & Model Interpretation

**03**

## **How much will it cost?**

- Cost Benefit Analysis

**04**

## **What should we do?**

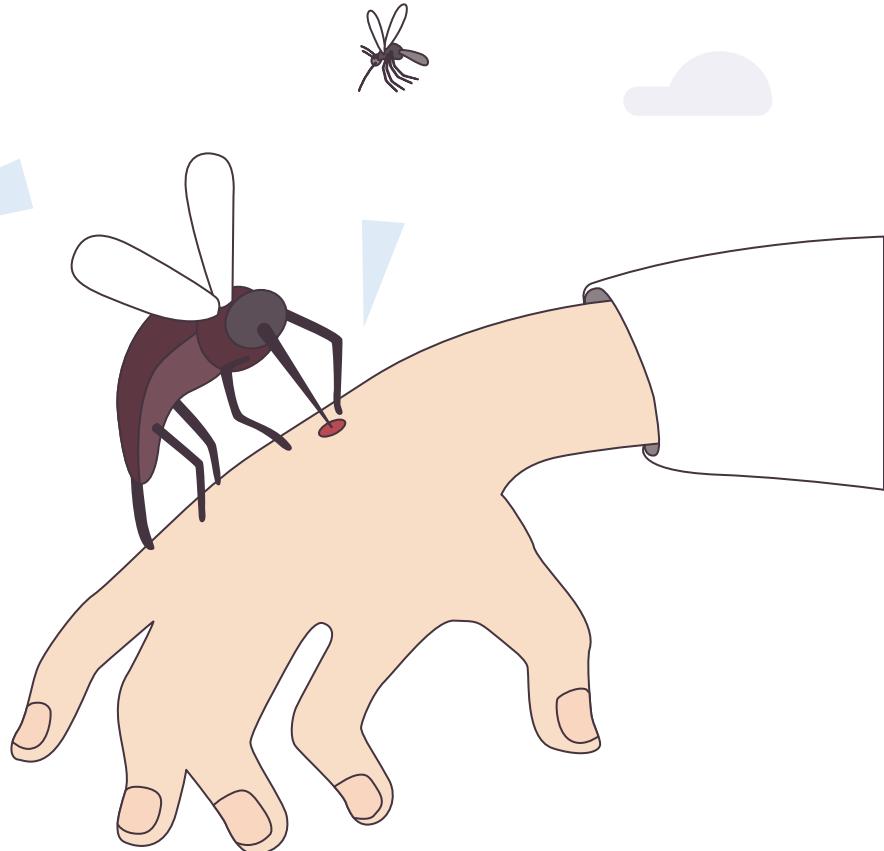
- Conclusions & Recommendations



# 01

# What's the big deal?

*Background & Motivations*



# 2746 cases / year

Across the United States of America<sup>1</sup>

## 1 in 5

Developed symptoms after being bitten

## \$56 million/year

In medical treatment and productivity losses<sup>2</sup>



South Dakota	GDP (millions \$): 55,243	% of USA: 0.30%	GDP Growth in 2018: 1.30%
--------------	---------------------------	-----------------	---------------------------

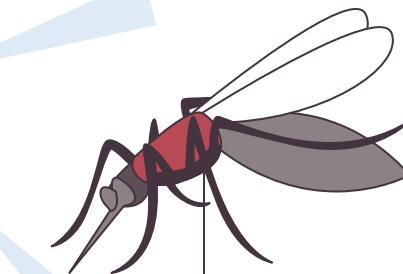
<sup>1</sup>

Average data from 2005 to 2020. Source: <https://www.cdc.gov/westnile/statsmaps/cumMapsData.html#three>

<sup>2</sup>

Average cost from 1999-2012. Source: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3945683/>

# Deadly Characteristics

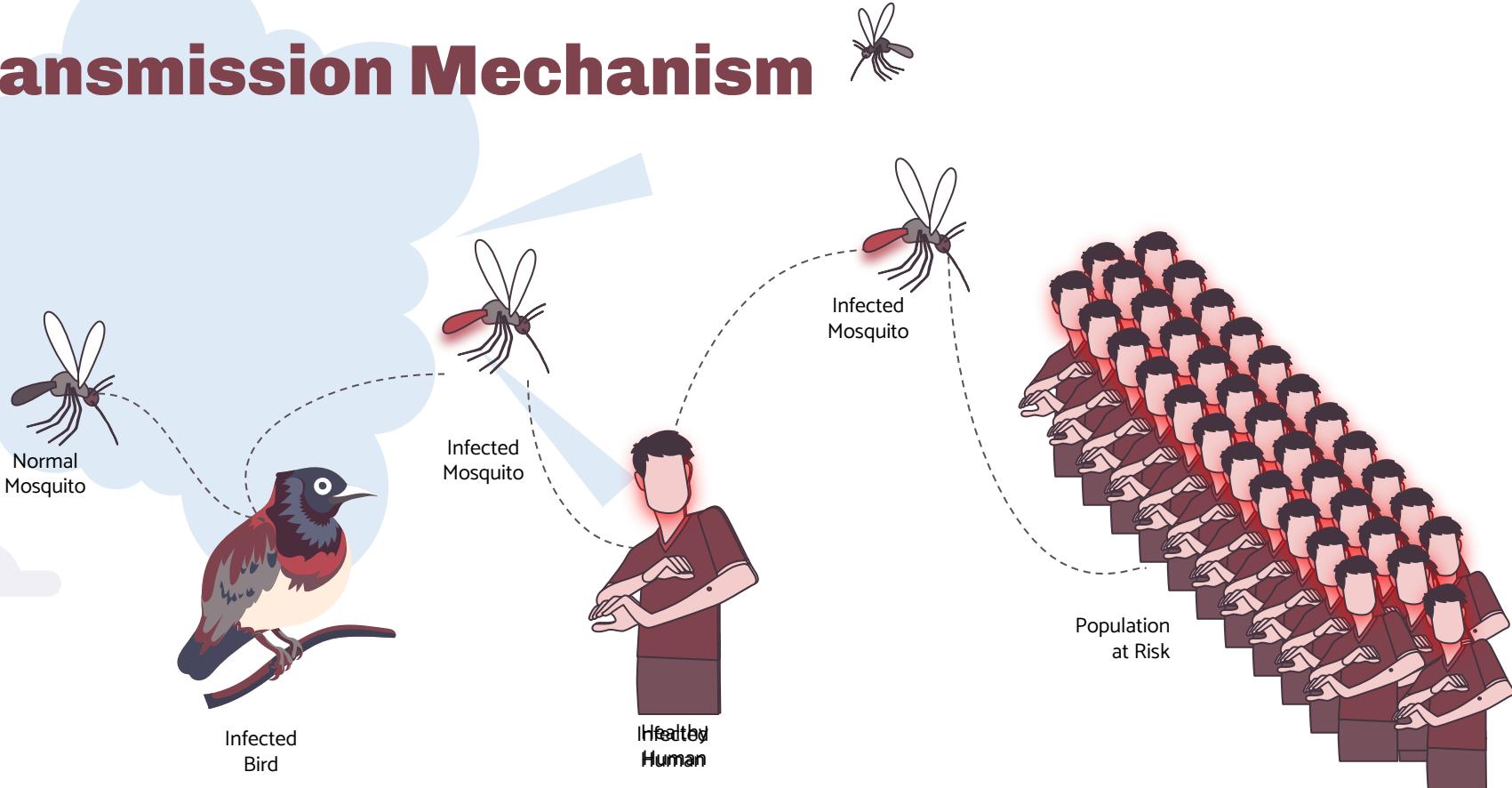


**Transmission  
Mechanism**

**High Mobility  
Carriers**

**Cost of  
Treatment**

# Transmission Mechanism





# **Project Objectives**

**Estimate the number of mosquitoes in an area**

**Predict the presence of WNV in an area**

**Estimate the costs of vector controls**

**Recommend Time and Place to conduct spraying**



# 02

## What do we know?

*Data Collection, & EDA*



# Data Provided



## Data Source

Chicago Department of  
Public Health



## Date Range

Spring - Summer  
2007 - 2014

**50-50 split** by time for  
training and validation



## Traps

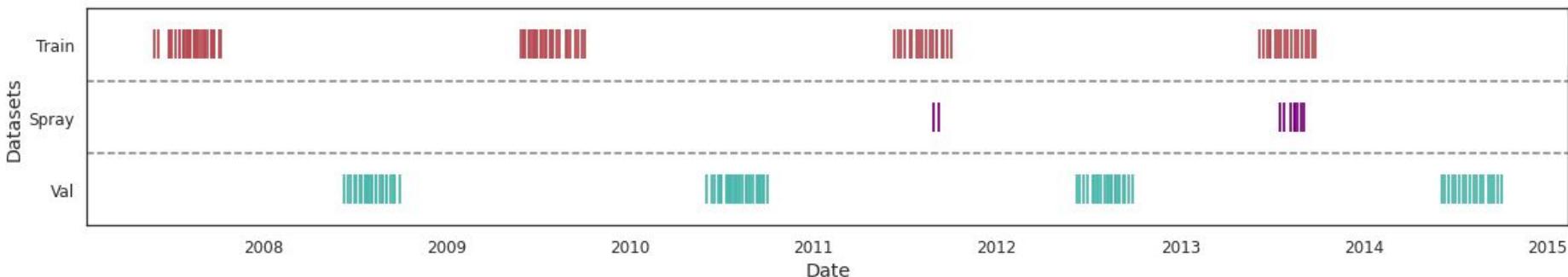
**136** Mosquito Traps  
Number of Mosquitoes  
WNV Presence  
Sprays Conducted



## Weather

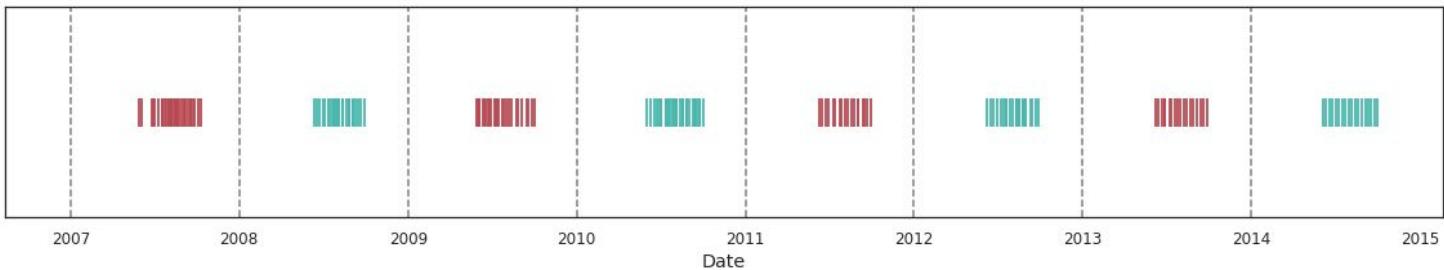
2 Weather Stations  
Temperature  
Wind Direction/Speed  
Humidity

# Date Range



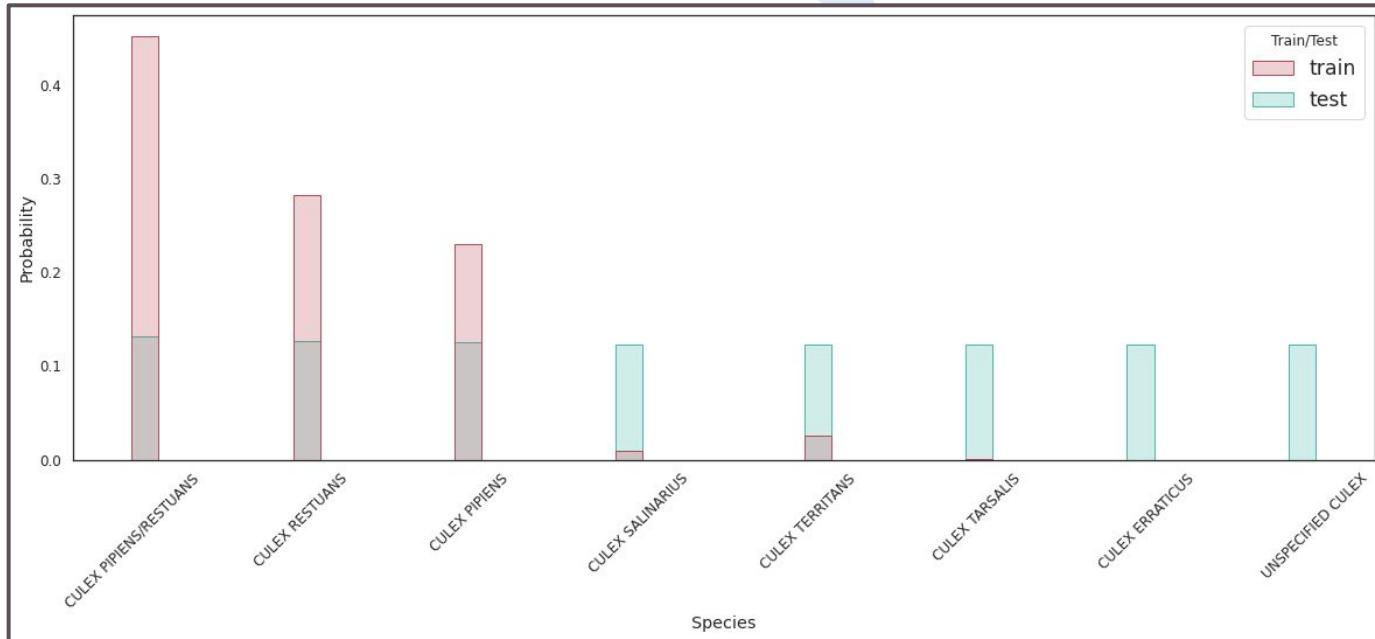
- Training and Validation data sets are split across alternate years
- Spray data was collected for only 2 of 4 years in the training data set

# Date Range



- Training and Validation data sets are split across alternate years
- Spray data was collected for only 2 of 4 years in the training data set

# Species Distribution



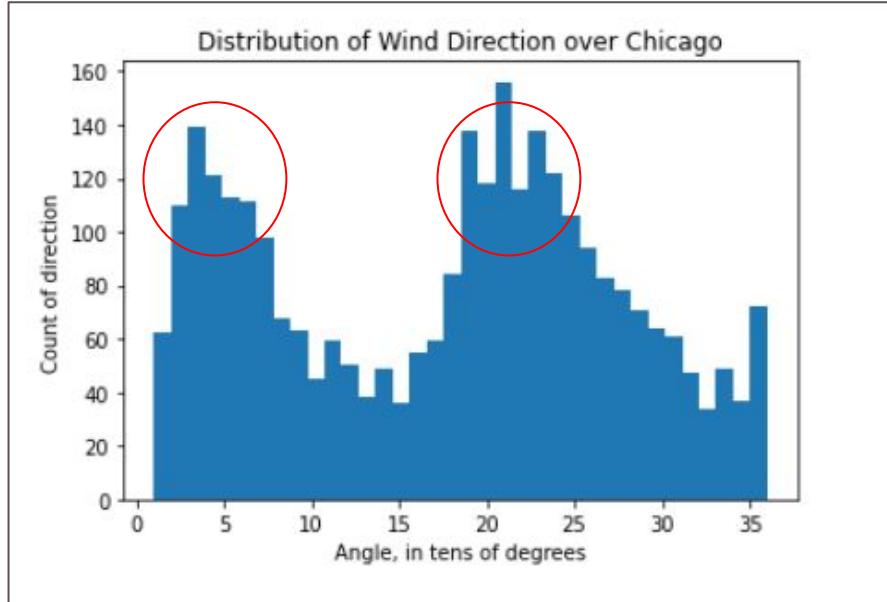
- CULEX PIPiens & CULEX RESTUANS account for more than 99.5% of the data set

# Species Distribution



- CULEX PIPiens & CULEX RESTUANS account for more than 99.5% of the data set

# Wind Direction over Chicago



- Data provided was rounded off to **nearest 10 degrees.**
- Prevailing winds are from predominantly from the **~20° and ~210° direction.**



Chicago

Illinois  
USA

Partly cloudy 7°C  
2:43 AM

Directions Save Nearby Send to your phone Share

**Quick facts**

Chicago, on Lake Michigan in Illinois, is among the largest cities in the U.S. Famous for its bold architecture, it has a skyline punctuated by skyscrapers such as the iconic John Hancock Center, 1,451-ft. Willis Tower (formerly the Sears Tower) and the neo-Gothic Tribune Tower. The city is also renowned for its museums, including the Art Institute of Chicago with its noted Impressionist and Post-Impressionist works.

**Iconic Chicago**

The Art Institute of Chicago  
4.8 ★ (26,292)  
Renowned art museum with global works

Navy Pier  
4.6 ★ (60,360)  
Destination with

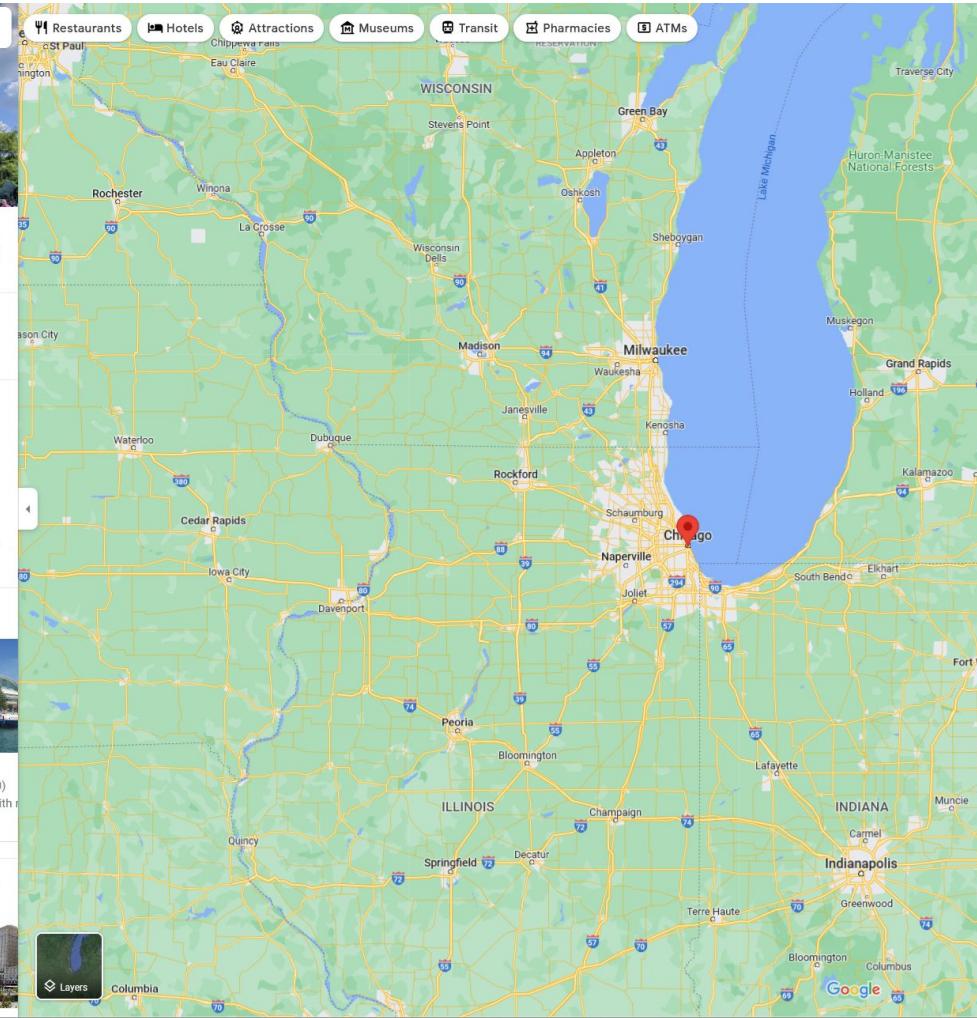
**Hotels**

About pricing ⓘ

\$205

\$202

Restaurants Hotels Attractions Museums Transit Pharmacies ATMs



Restaurants Hotels

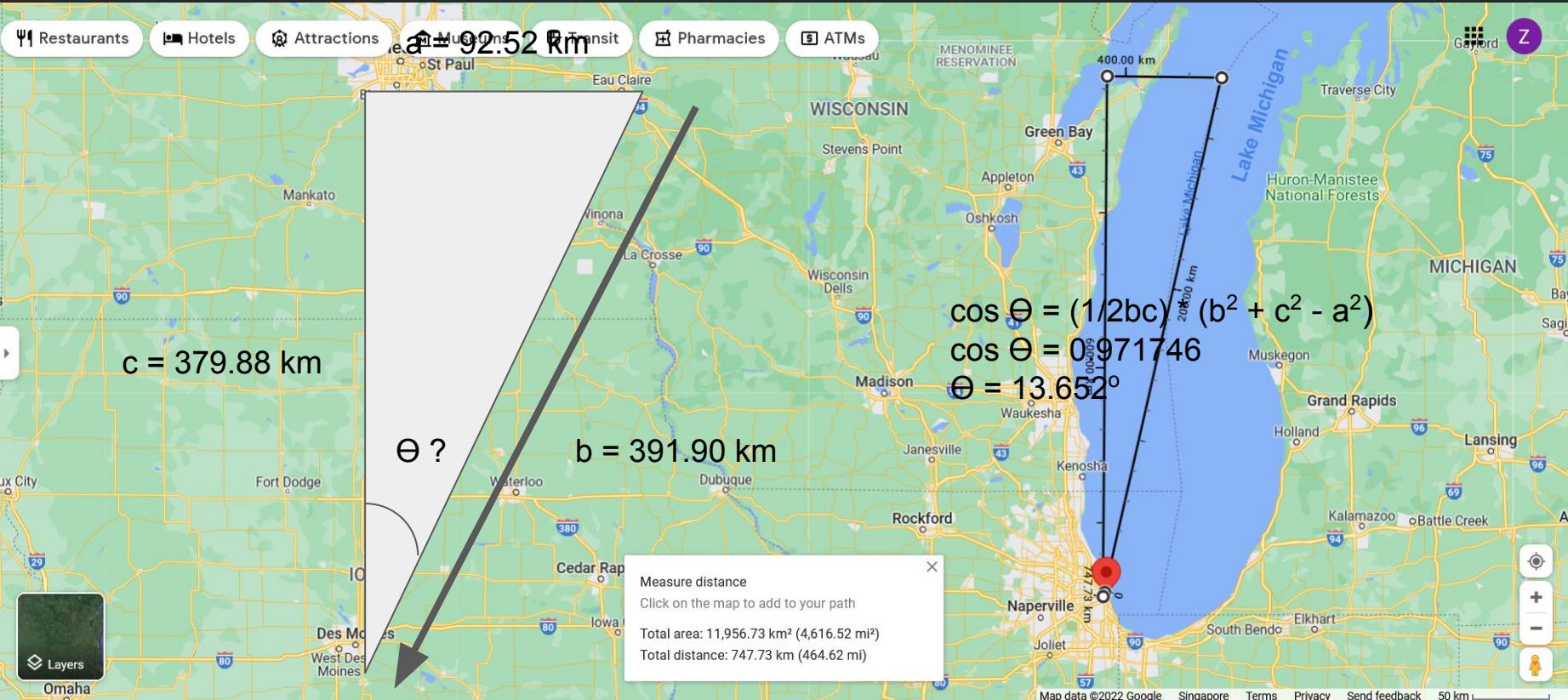
Attractions Museums

Transit  $a = 92.52 \text{ km}$

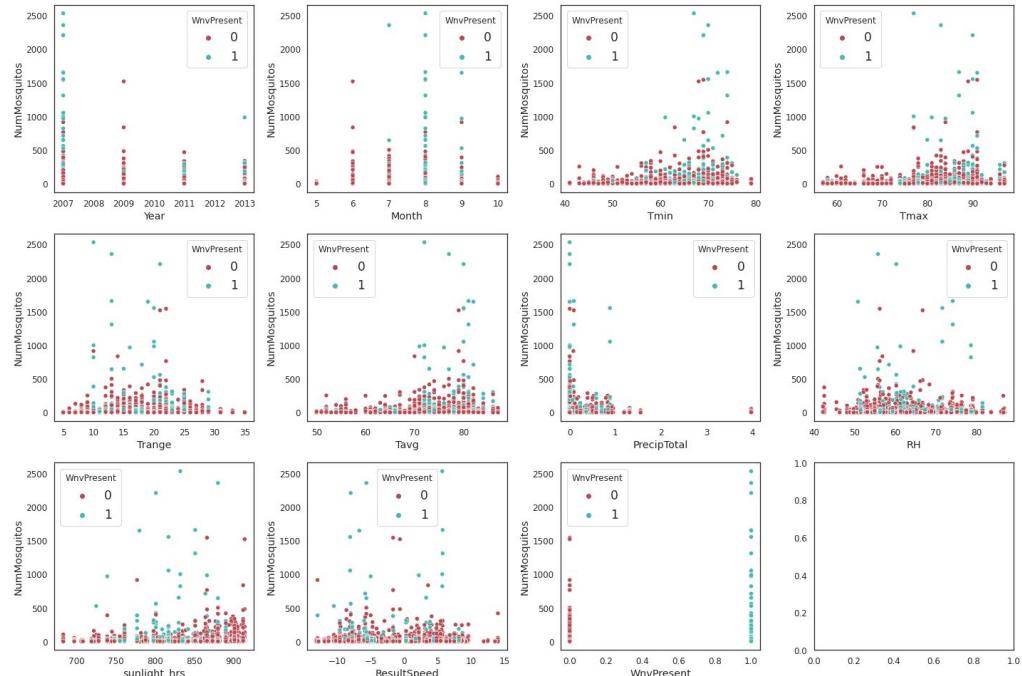
Pharmacies

ATMs

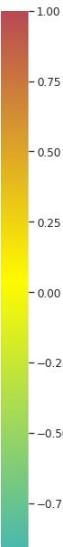
Z



# No Strong Linear Correlation with Response



	Year	Month	Tmin	Tmax	Trange	Tavg	PrecipTotal	RH	sunlight_hrs	ResultSpeed	NumMosquitos	WnvPresent
Year	1	-0.14	-0.07	-0.062	0.0043	-0.07	-0.087	-0.13	0.13	0.11	-0.02	0.043
Month	-0.14	1	-0.076	-0.013	0.091	-0.045	-0.14	-0.012	-0.92	-0.19	-0.0082	0.098
Tmin	-0.07	-0.076	1	0.79	-0.2	0.94	0.17	0.065	0.22	-0.32	0.064	0.08
Tmax	-0.062	-0.013	0.79	1	0.44	0.95	0.049	-0.22	0.16	-0.49	0.055	0.059
Trange	0.0043	0.091	-0.2	0.44	1	0.14	-0.18	-0.45	-0.074	-0.32	-0.0062	-0.024
Tavg	-0.07	-0.045	0.94	0.95	0.14	1	0.12	-0.095	0.2	-0.43	0.063	0.073
PrecipTotal	-0.087	-0.14	0.17	0.049	-0.18	0.12	1	0.46	0.14	-0.15	-0.0045	0.014
RH	-0.13	-0.012	0.065	-0.22	-0.45	-0.095	0.46	1	-0.053	0.0042	-0.0065	0.051
sunlight_hrs	0.13	-0.92	0.22	0.16	-0.074	0.2	0.14	-0.053	1	0.17	0.026	-0.073
ResultSpeed	0.11	-0.19	-0.32	-0.49	-0.32	-0.43	-0.15	0.0042	0.17	1	-0.0038	-0.014
NumMosquitos	-0.02	-0.0082	0.064	0.055	-0.0062	0.063	-0.0045	-0.0065	0.026	-0.0038	1	0.23
WnvPresent	0.043	0.098	0.08	0.059	-0.024	0.073	0.014	0.051	-0.073	-0.014	0.23	1

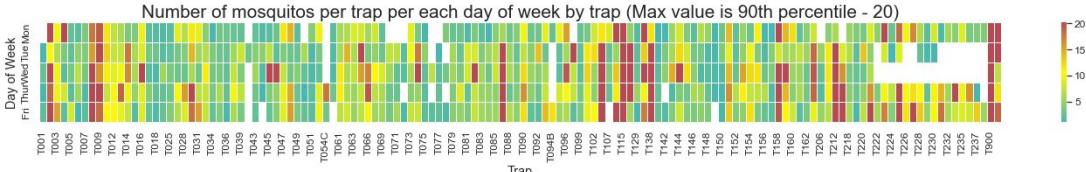
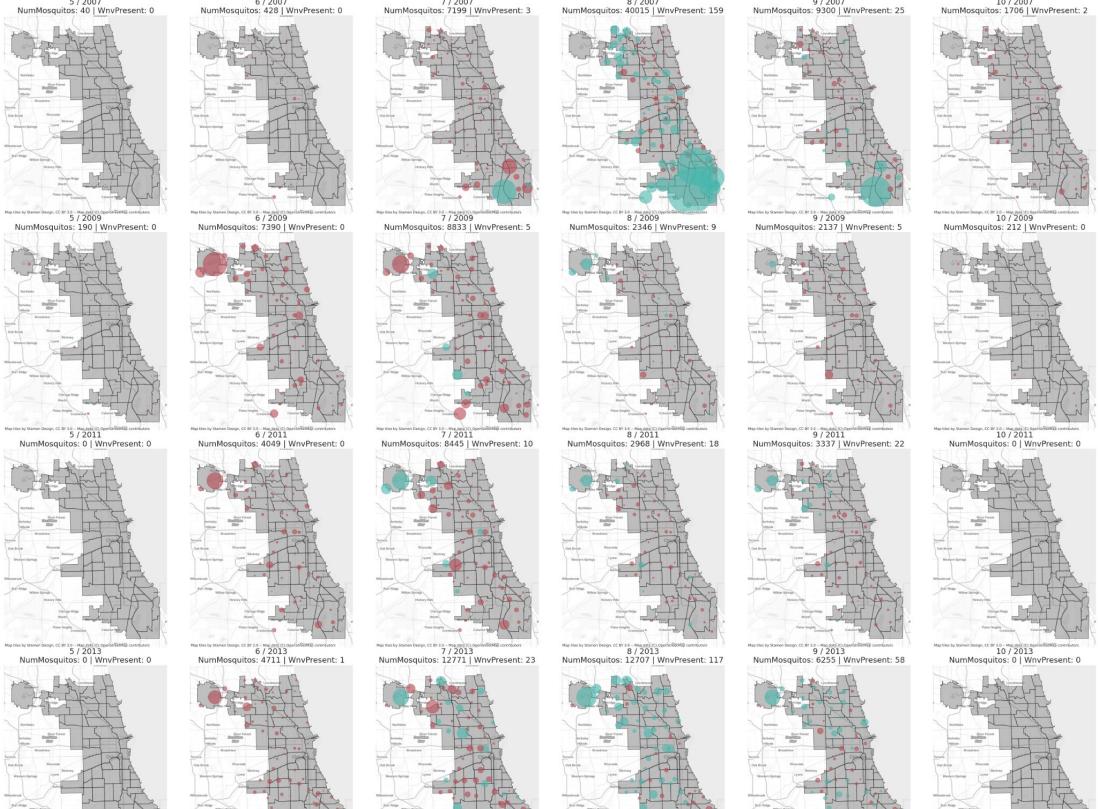


# Hotspots

- Hotspots tend to result in rapid growth rate of the number of mosquitoes
- The high number of mosquitoes puts the area at a higher risk of WNV prevalence



Size: Higher  
NumMosquitos



# Hotspots

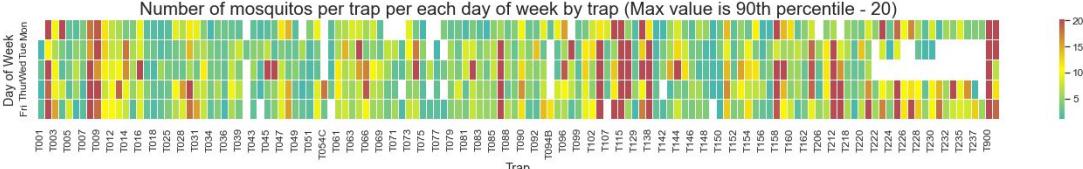
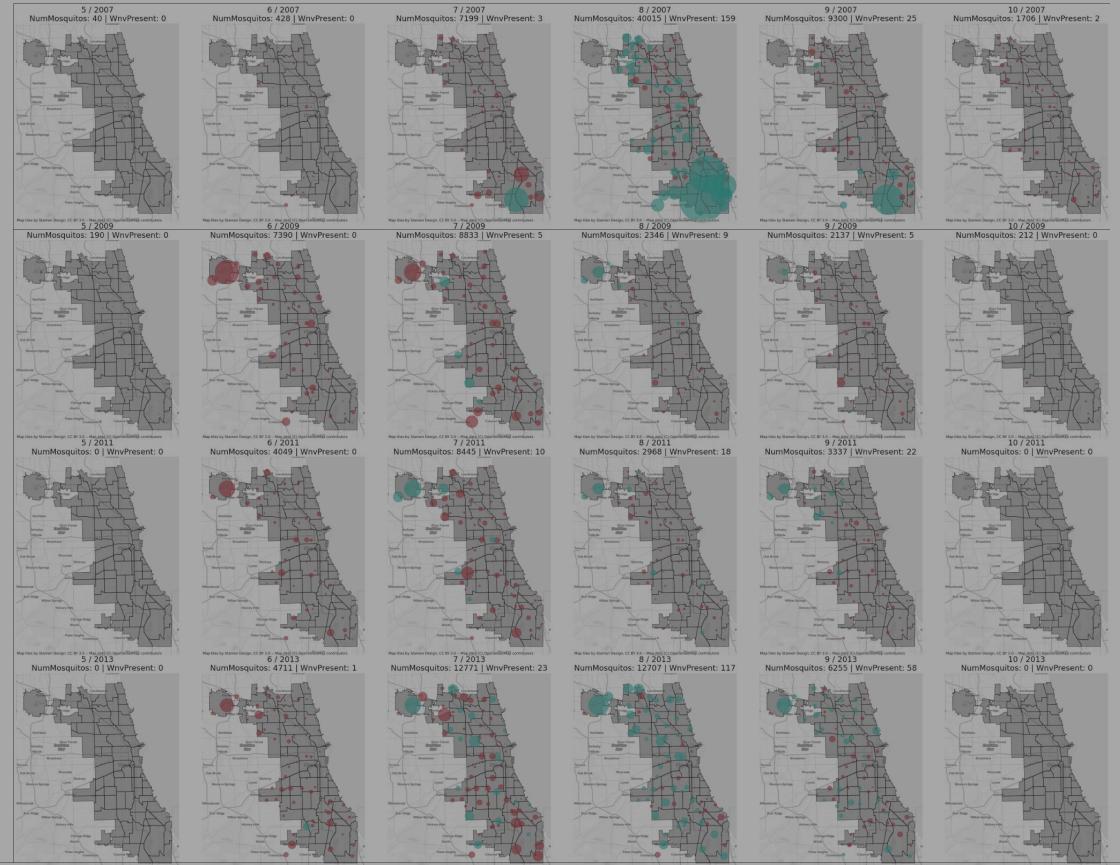
- Hotspots tend to result in rapid growth rate of the number of mosquitoes
- The high number of mosquitoes puts the area at a higher risk of WNV prevalence



Size: Higher  
NumMosquitos



Wnv  
Present



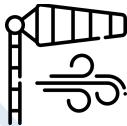


# 02

## What do we know?

*Feature Engineering & Data Modelling*

# Feature Engineering

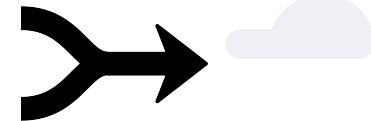


## Wind Direction

Resolving wind direction towards **Lake Michigan**

## One-Hot-Encode

One hot encoding of categorical features



## Species Aggregation

Aggregating based on  
*Culex Pipen/Restuans*



## Lag Weather

Lagging weather observations by **5 to 14** days

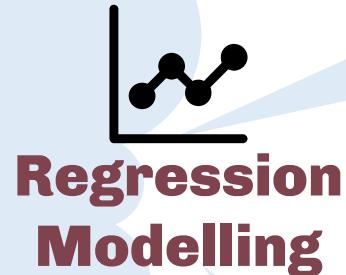
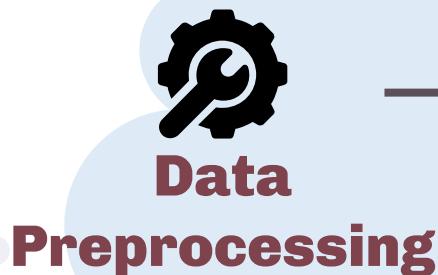


## Distance Matrix

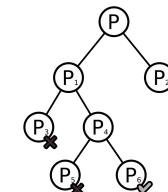
Inverse weighted distance from neighbouring traps



# Modelling Approach



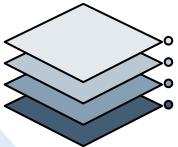
To obtain  
NumMosquitos



To determine  
Wnv Presence

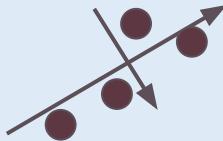


# Preprocessing



## Species Balancing

Add rows for unfound species



## PCA

Conduct PCA for linear models



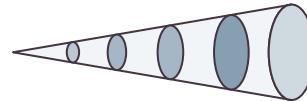
## Oversampling

Oversampling on minority response class using SMOTE



## Cross validation

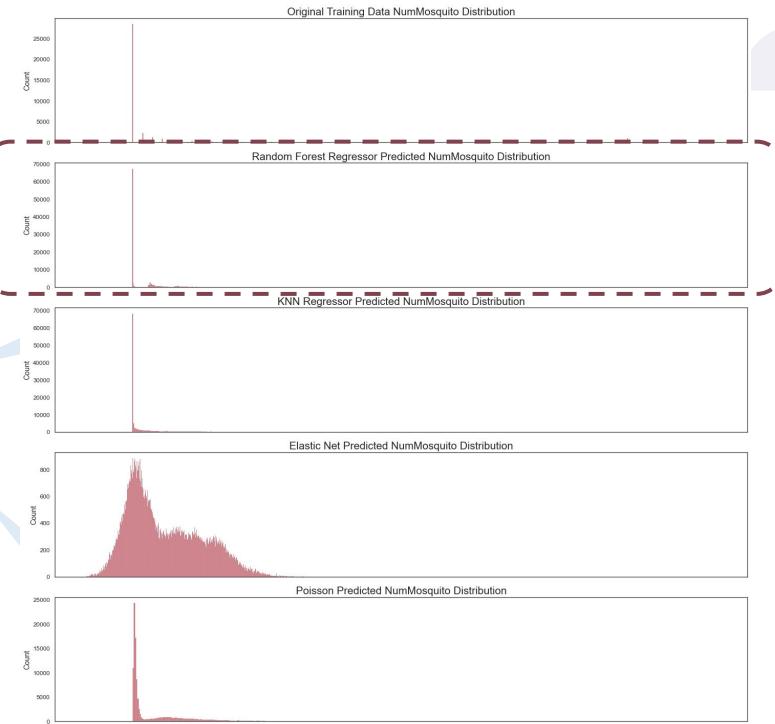
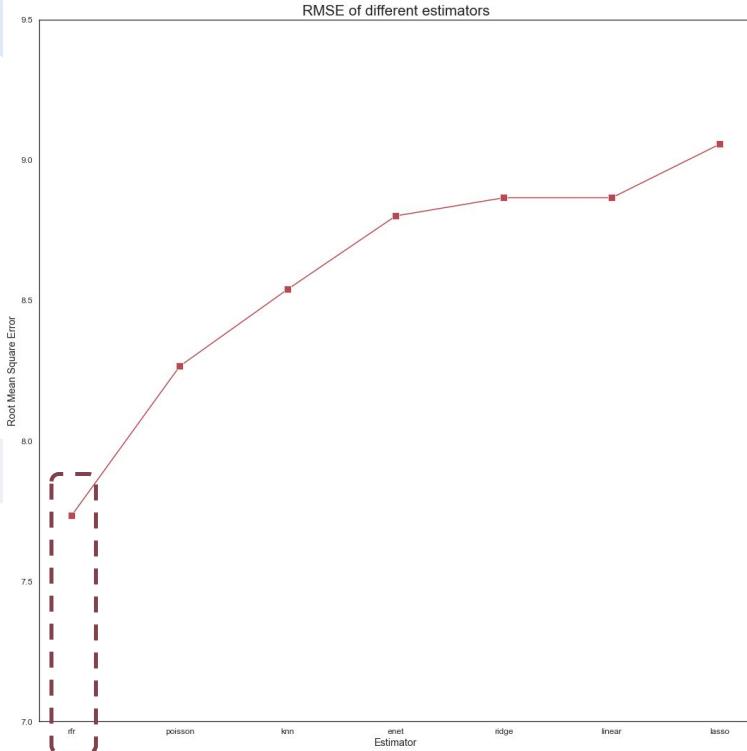
5-fold Cross-validation on training set



## Scaling

Either Standard Scaling or MinMax Scaling

# Model Performance (Regression)



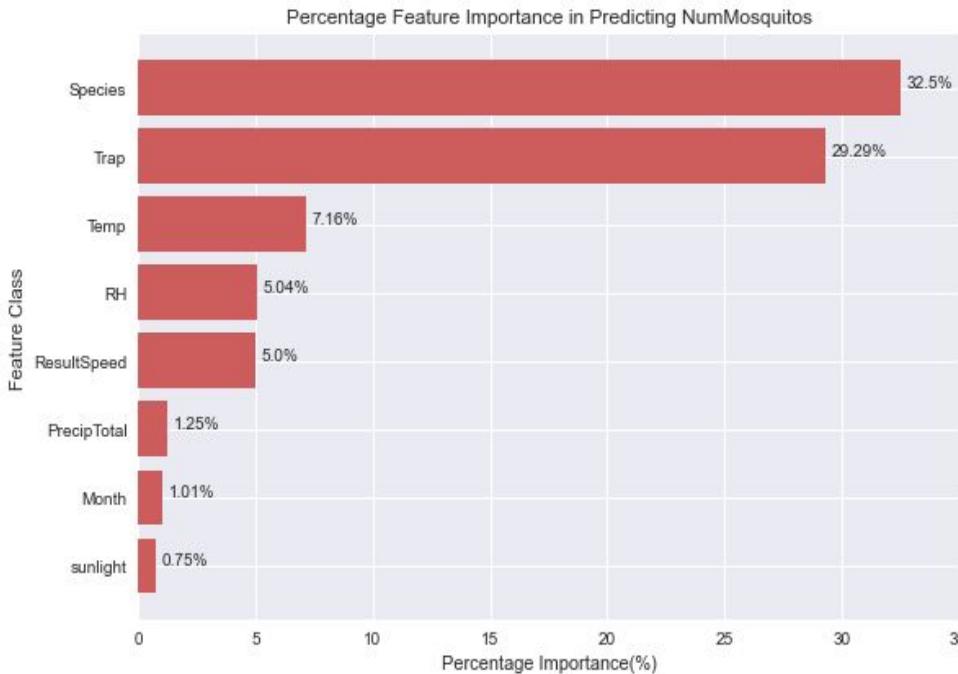
**Chosen Model: Random Forest Regressor**

# Model Performance (Classification)

Cross Validated Scores

Model	ROC-AUC	Accuracy	Recall	Precision
Gradient Boosting	0.93	0.93	0.55	0.13
Random Forest	0.91	0.93	0.70	0.14
Logistics Regression	0.77	0.88	0.52	0.07

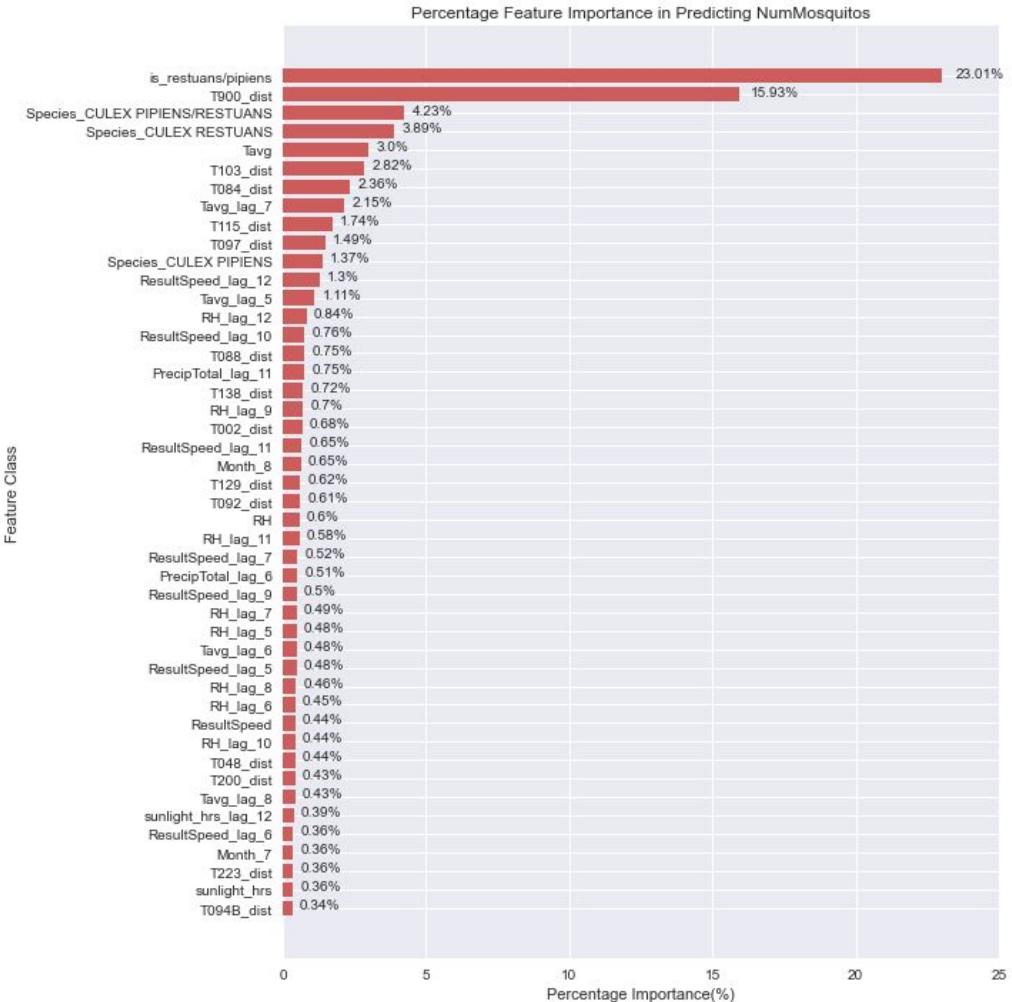
# Feature Importance: NumMosquitos



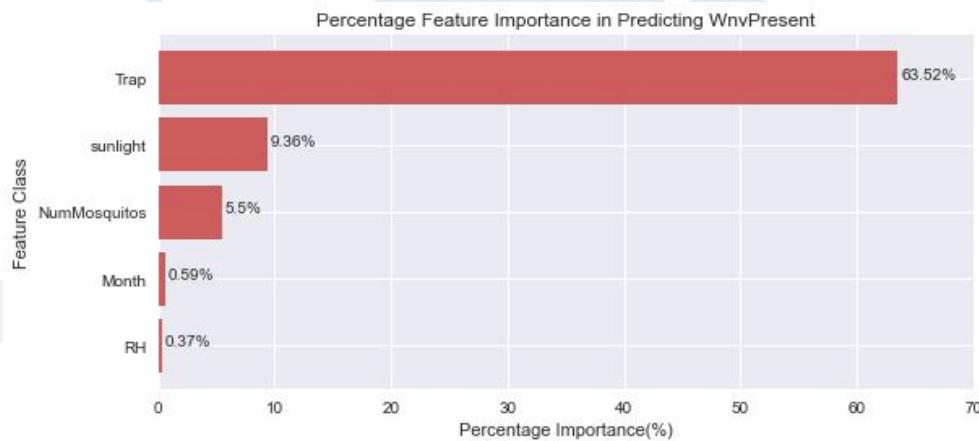
- Species type (32.5%)
  - Pipiens / Restuans
- Trap location (29.3%)
- Weather Data (18.5%)
  - Temperature
  - Humidity
  - Wind Speed/Direction
  - Total Rainfall
- Three categories sum to >80%



# Feature Importance: NumMosquitos



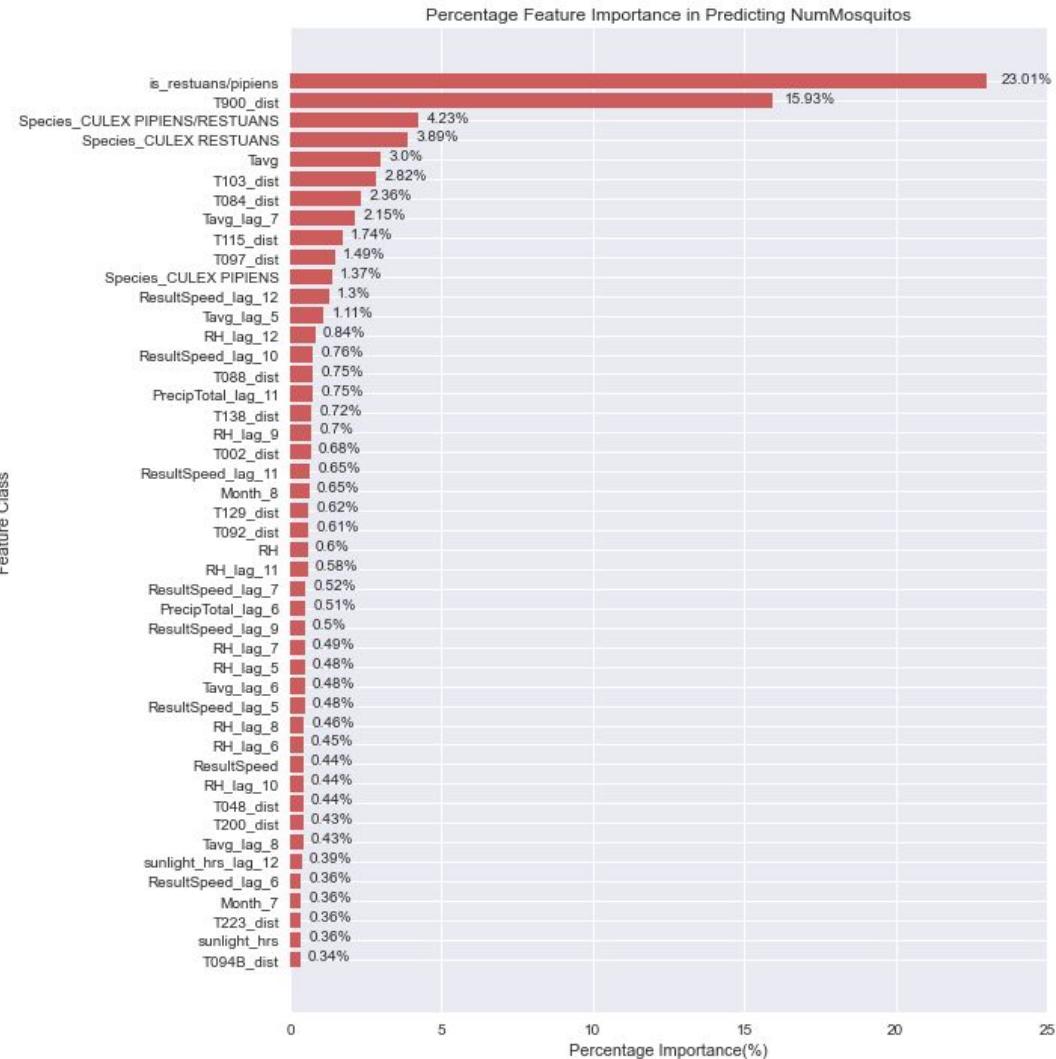
# Feature Importance: WnvPresent



- Trap location (63.52%)
- Sunlight hours (9.36%)
- NumMosquitos (5.5%)
- Month Number (0.59%)



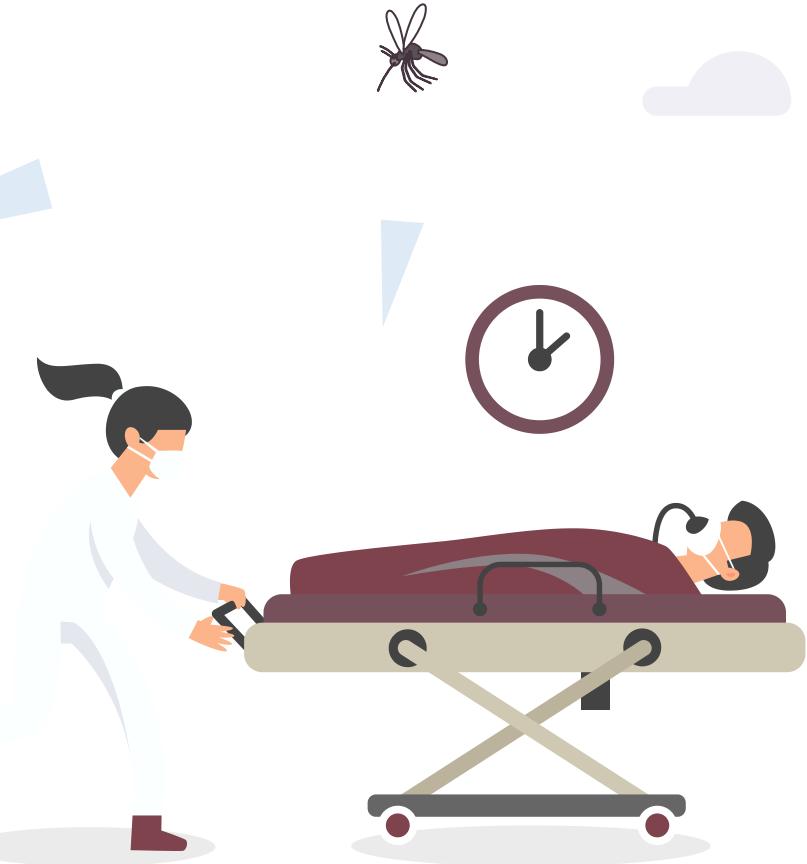
# Feature Importance: WnvPresent



# 03

# How much will it cost?

*Cost Benefit Analysis*



# Primary Cost Drivers



## Spraying

Cost associated with conducting vector controls



## Healthcare & Productivity

Costs associated with medical coverage & productivity losses



# Spray Cost Model



## Spray Costs

Cost associated with conducting vector controls

<b>f</b>	<b>Frequency</b>	Number of sprays in a given period
<b>A</b>	<b>Area</b>	Area covered
<b>UC<sub>s</sub></b>	<b>Unit Cost</b>	Cost per spray

$$C_s = f \times A \times UC_s$$



# Healthcare Cost Model



## Healthcare & Productivity Costs

Costs associated with medical coverage & productivity losses

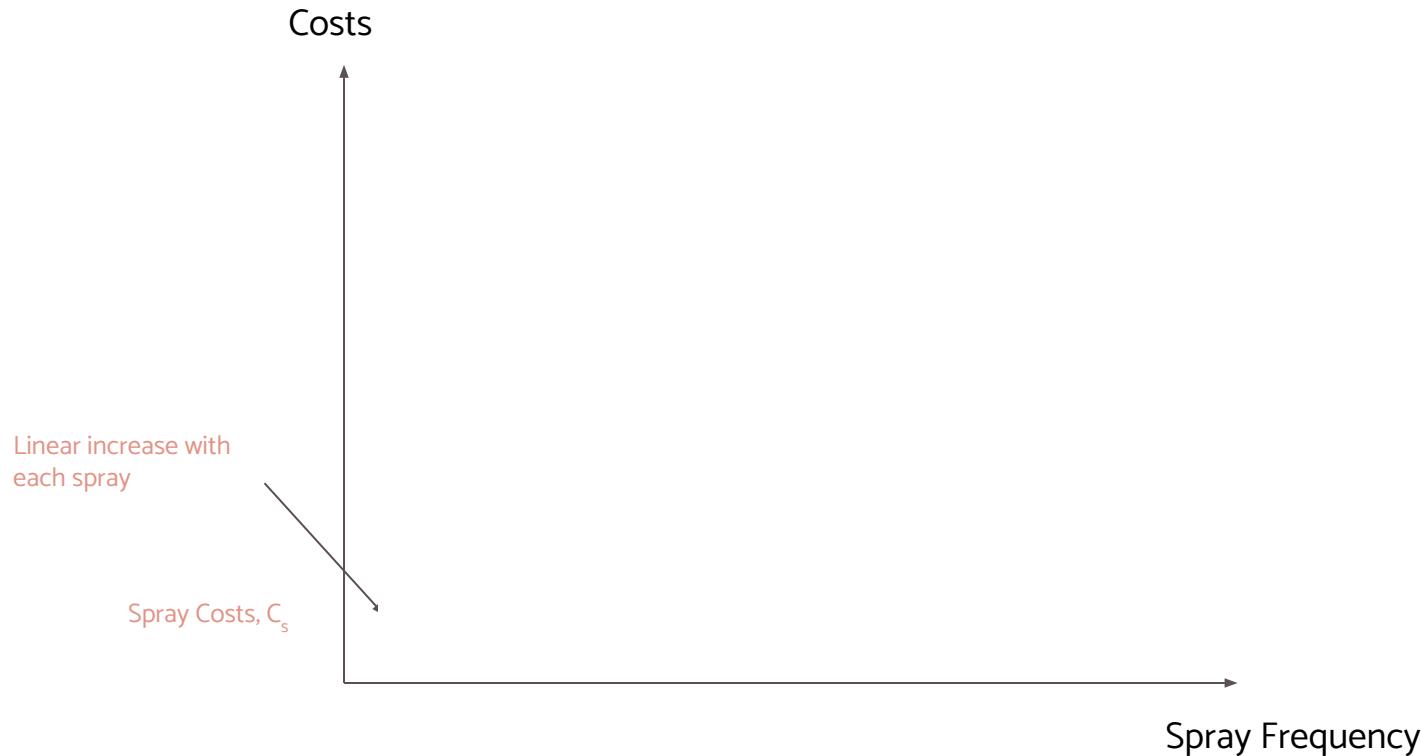
<b>N<sub>M</sub></b>	<b>NumMos</b>	Number of mosquitoes trapped
<b>T</b>	<b>%Trapped</b>	Percentage of mosquitoes trapped
<b>N<sub>B</sub></b>	<b>NumBites</b>	Number of human bites, per mosquito
<b>V</b>	<b>WnvPresent</b>	Binary toggle; Indicates if WNV is detected in mosquitoes trapped
<b>UC<sub>H</sub></b>	<b>Unit Cost</b>	Averaged healthcare and productivity costs, per person

$$C_H = (N_M \times T) \times (N_B \times V) \times UC_H$$

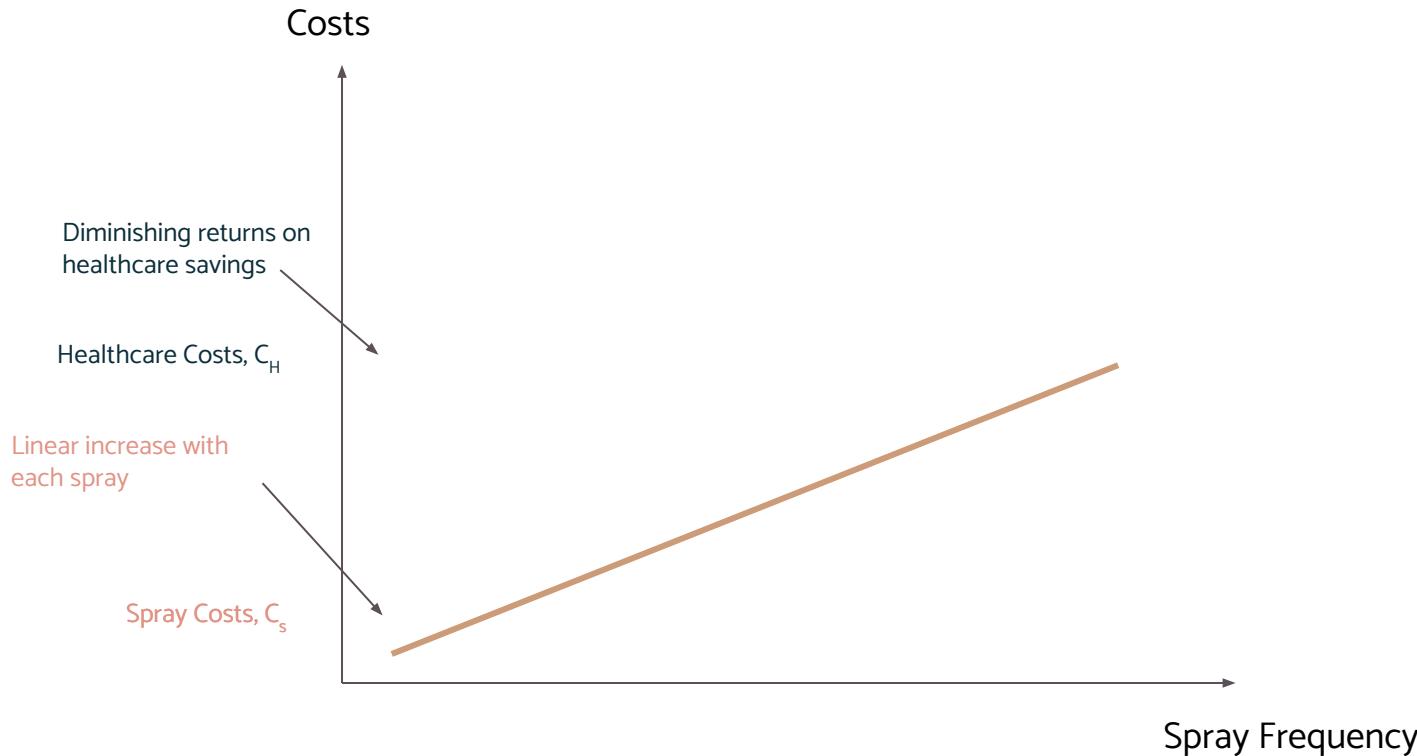
$$C_H = \text{Mosquito Population} \times \text{Human Infection Rate} \times \text{Per Head Cost}$$



# Cost Optimization

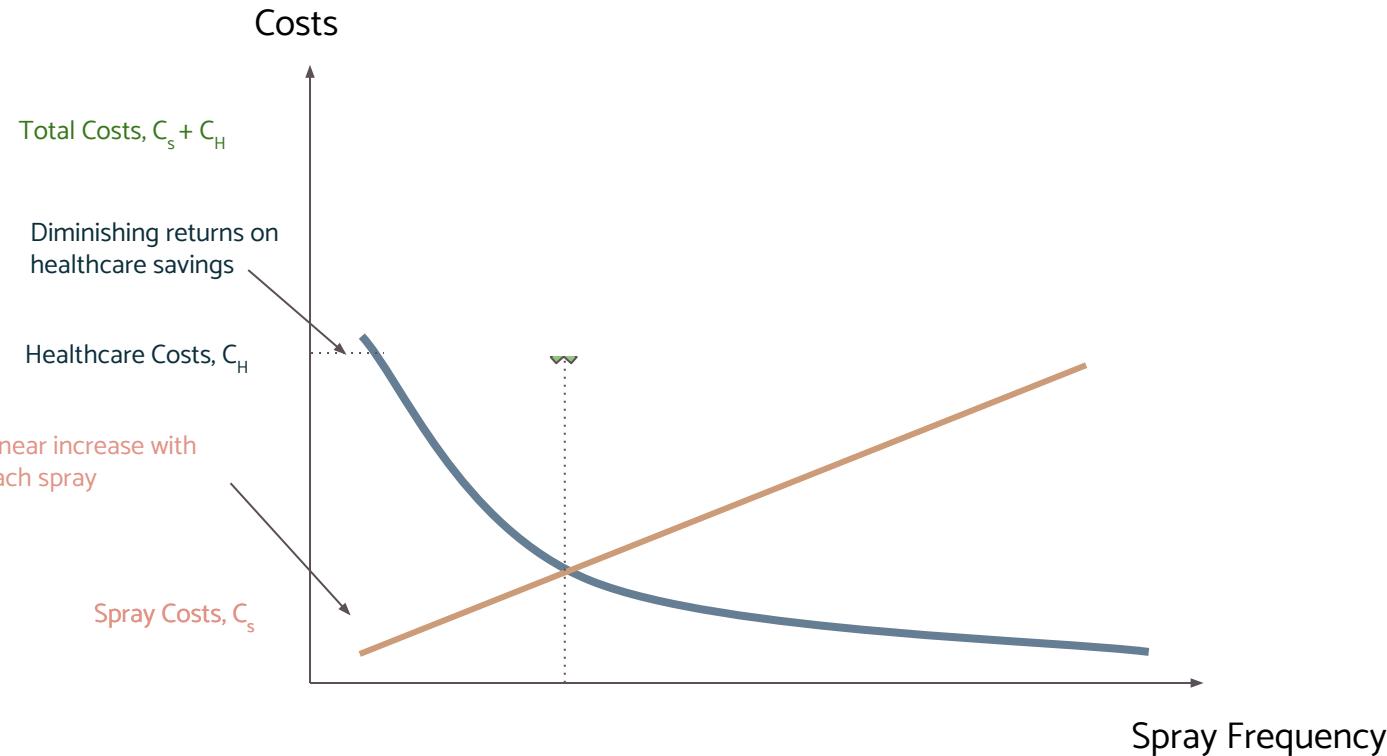


# Cost Optimization

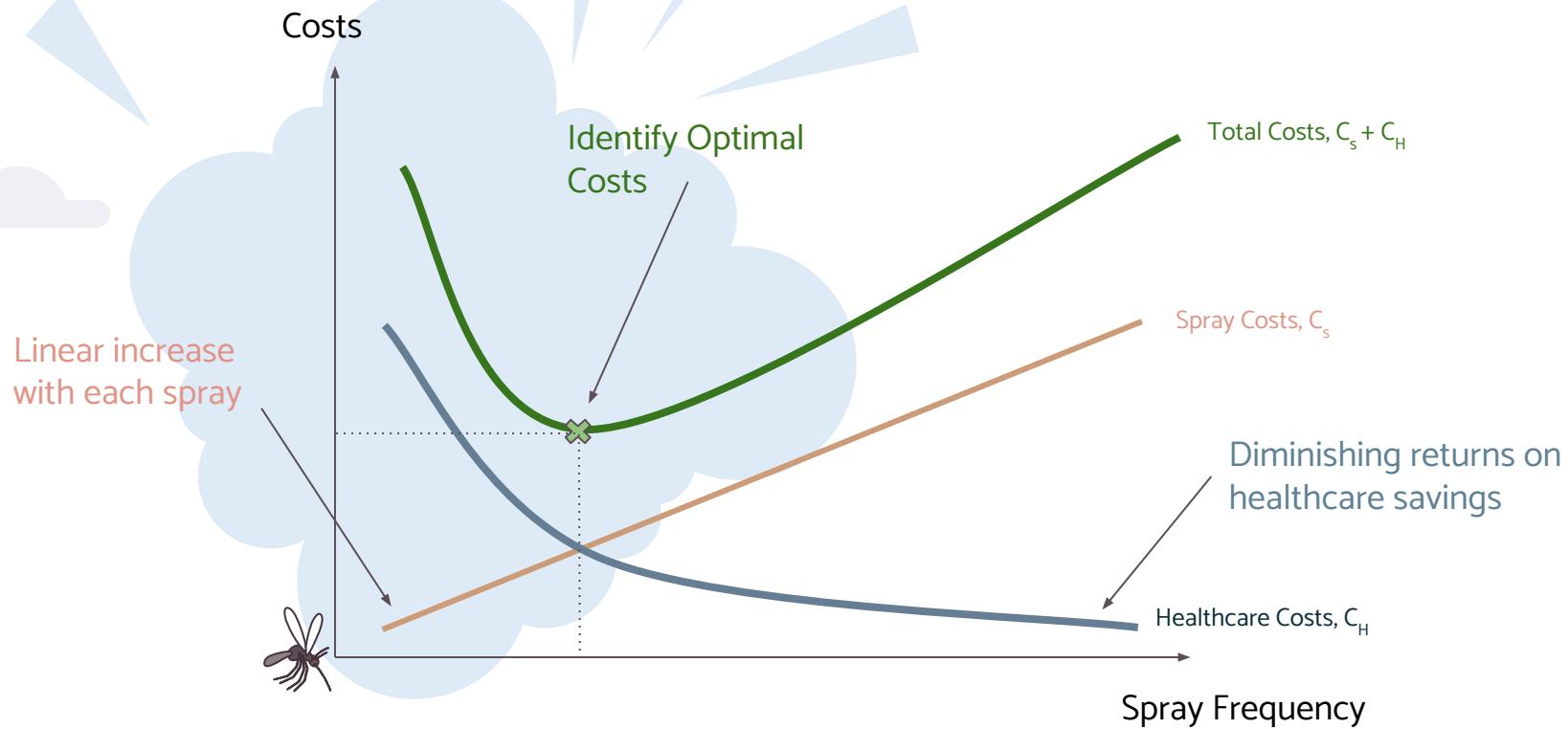


# Cost Optimization

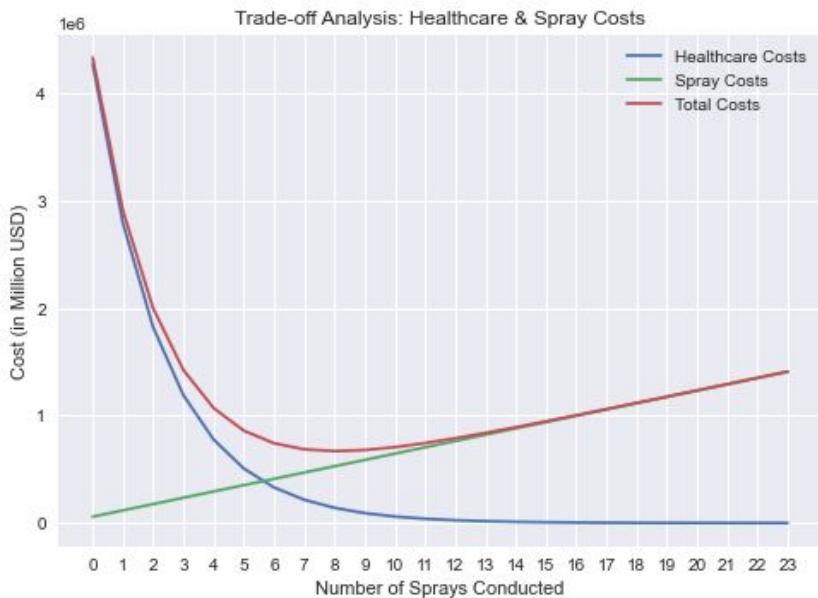
Identify Optimal Costs



# Cost Optimization



# Simulation Results



**Min Total Costs**

\$669,321

**Optimal Spray Freq**

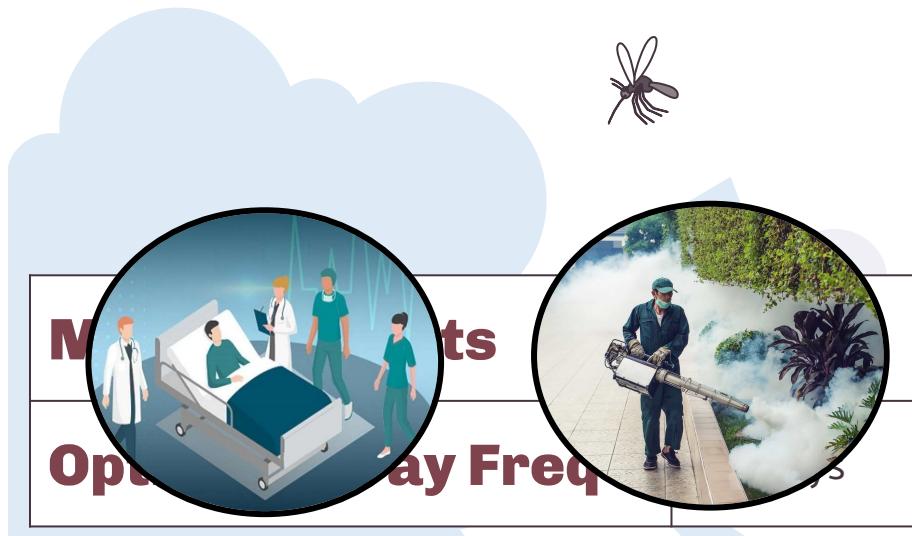
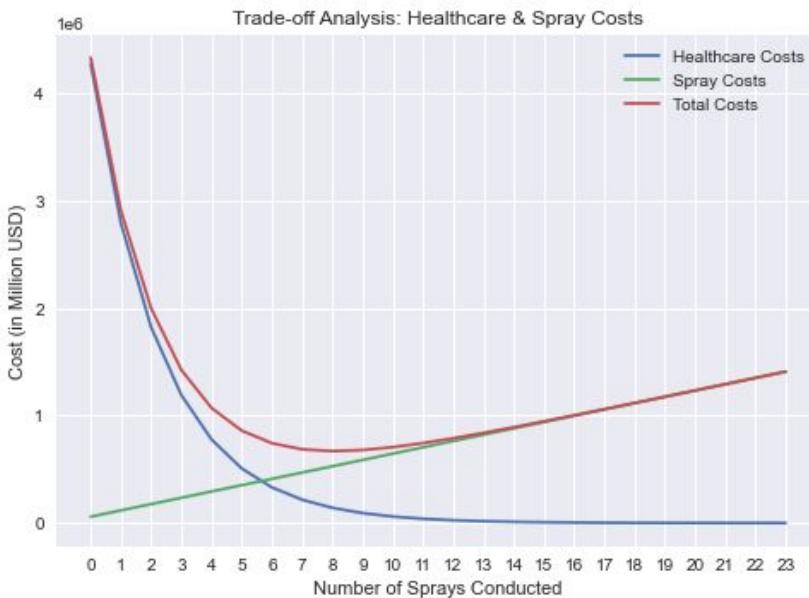
9 sprays



Healthcare Costs:  
%Population Infected: \$6,317 USD per person  
0.1% (1,000 per million)

Spray Costs:  
%Infected Reduction: \$0.92 / acre  
65.3%

# Simulation Results



# 04

## What should we do?

*Conclusions & Recommendations*



# Conclusions



## High Risk in July-Sept

Consistently High Wnv Occurrence in these months



## Hotspots Identified

Traps were ranked based on relative importance



# Recommendations

## Preventive Sprays

Aggressive spraying in earlier months

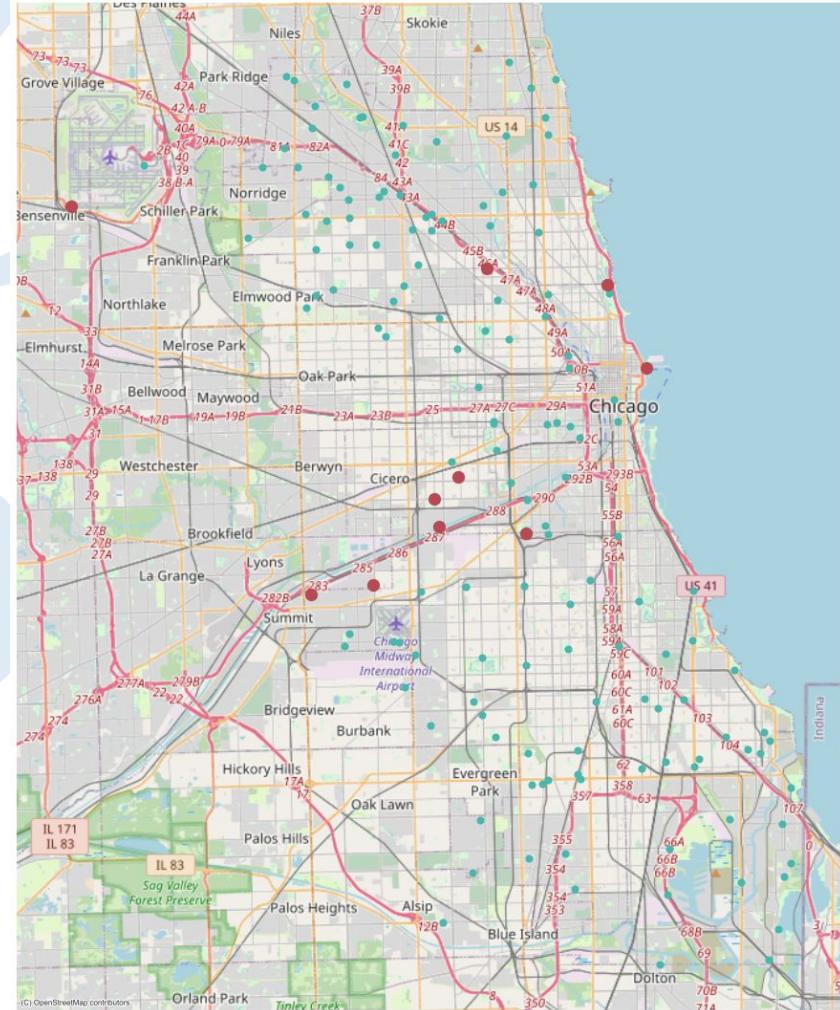


## Increased Surveillance

Increasing surveillance near hotspots



# **Location of Critical Monitoring Traps**



# Thanks

Do you have any questions?



# Future Improvements

- From the plots, we visually see that the spray areas are concentrated in distinct clusters
  - Incorporate Local Population as inputs
    - Has a key role in outbreak mechanisms
  - Incorporate Bird Population & Migratory Paths
    - Leverage robust population modelling techniques

### Logistic growth equation:

$$\frac{dN}{dt} = rN \left(1 - \frac{N}{K}\right)$$

### Competitive Lotka–Volterra equations:

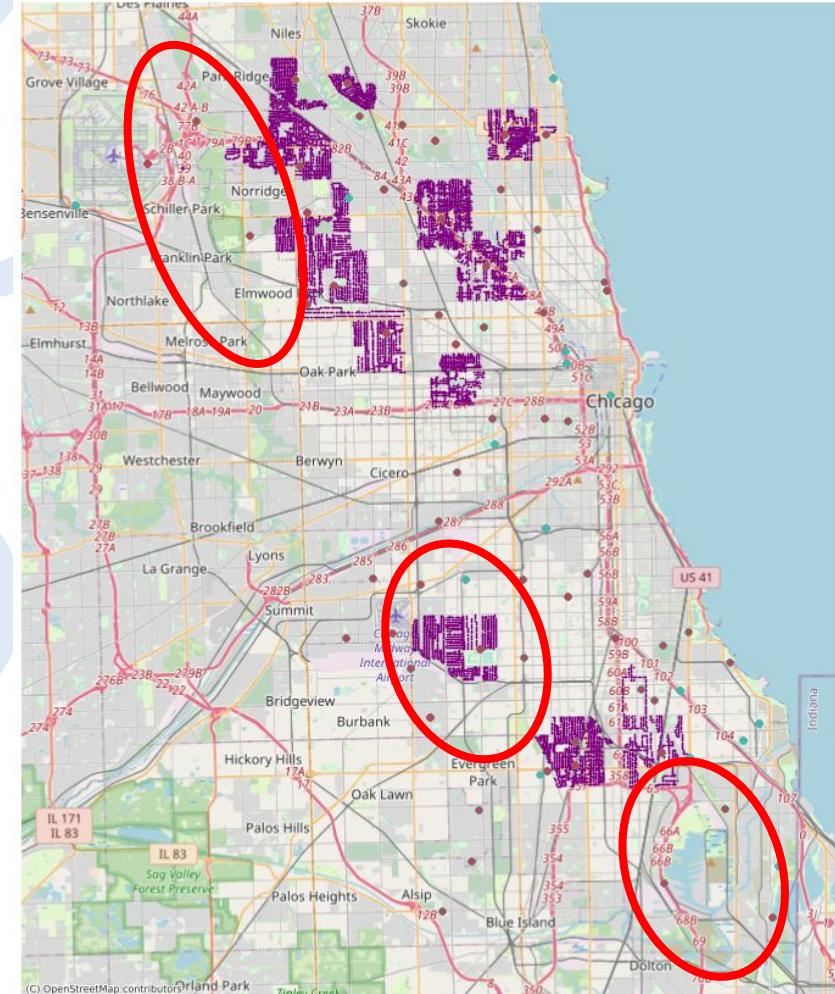
$$\frac{dN_1}{dt} = r_1 N_1 \frac{K_1 - N_1 - \alpha N_2}{K_1}$$

## Island biogeography:

$$S = \frac{IP}{I+E}$$

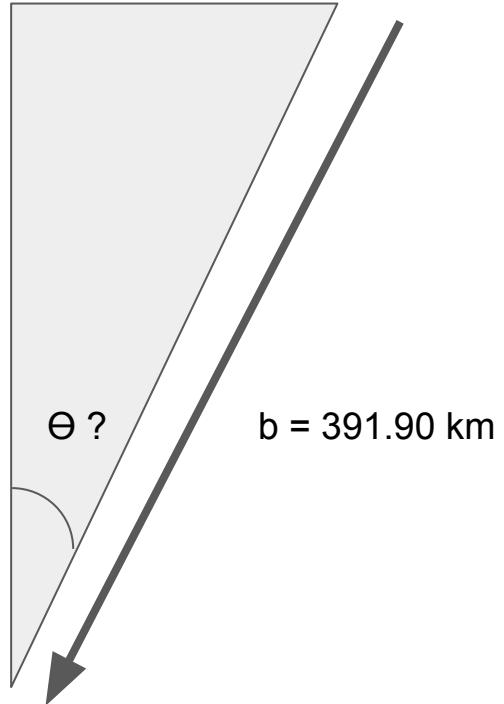
### Species-area relationship:

$$\log(S) = \log(c) + z \log(A)$$



$$a = 92.52 \text{ km}$$

$$c = 379.88 \text{ km}$$



$$\cos \Theta = (1/2bc) * (b^2 + c^2 - a^2)$$

$$\cos \Theta = 0.971746$$

$$\Theta = 13.652^\circ$$

# Fonts & colors used

This presentation has been made using the following fonts:

## Chivo

(<https://fonts.google.com/specimen/Chivo>)

## Catamaran

(<https://fonts.google.com/specimen/Catamaran>)

#443440

#5c4f58

#5d3742

#8c8184

#76515b

#7f434e

#b84a54

#cfe2f3

#ffffff