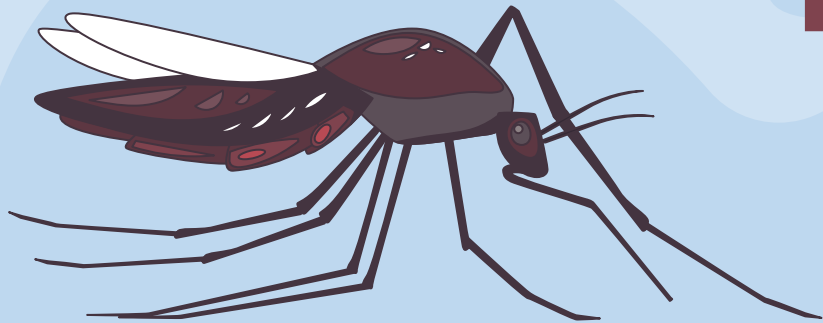


# **WNV** **Prediction** **in Chicago**



Emma  
Junhao  
Benjamin  
Gabriel

# Overview

- **Problem Statement**
- **Interesting Features**
- **Model Selection**
- **Feature Importance**
- **Cost Benefit Analysis**
- **Conclusion/Recommendations**

# Problem Statement

- We are a team of data scientists at the Disease And Treatment Agency. In the last few years, there have been several episodes of the West Nile Virus (WNV) outbreak in Chicago. The aim of our project is to leverage on weather, location and time data to predict West Nile Virus occurrences in Chicago.
- By **understanding these features and whether they can predict the virus**, we will then evaluate **whether spraying areas with pesticides will be an effective tool vis-a-vis medical and economic costs**

# Interesting Features



## Relative Humidity

Humidity - ambiguous relationship with WNV  
Instead of using DewPoint x Tavg interaction  
term - worked out  $RH=100 \times (e / e_s)$

## Lagged Weather Features

Culex: egg to adult ~7 to 10 days.  
Lagged weather features by 14 days - factor in  
time needed to pick up the virus



## Risk Classification

136 unique trap addresses  
5 categories based on WNV occurrences:  
Low, Very Low, Medium, High, Very High



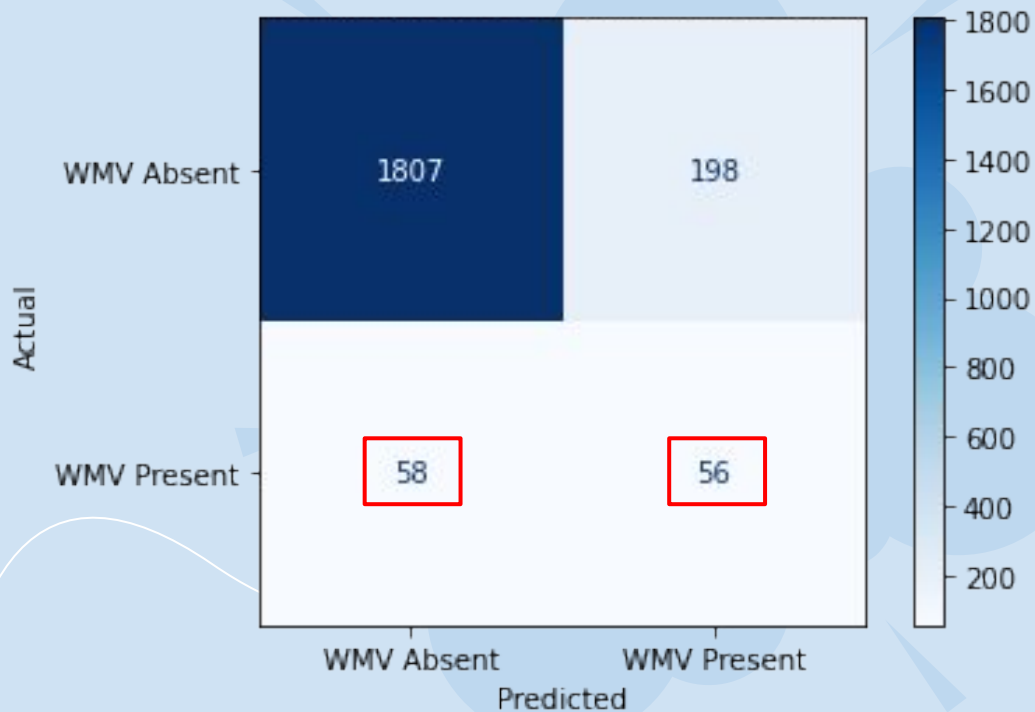
# Which model is suitable?

(Selected Model)

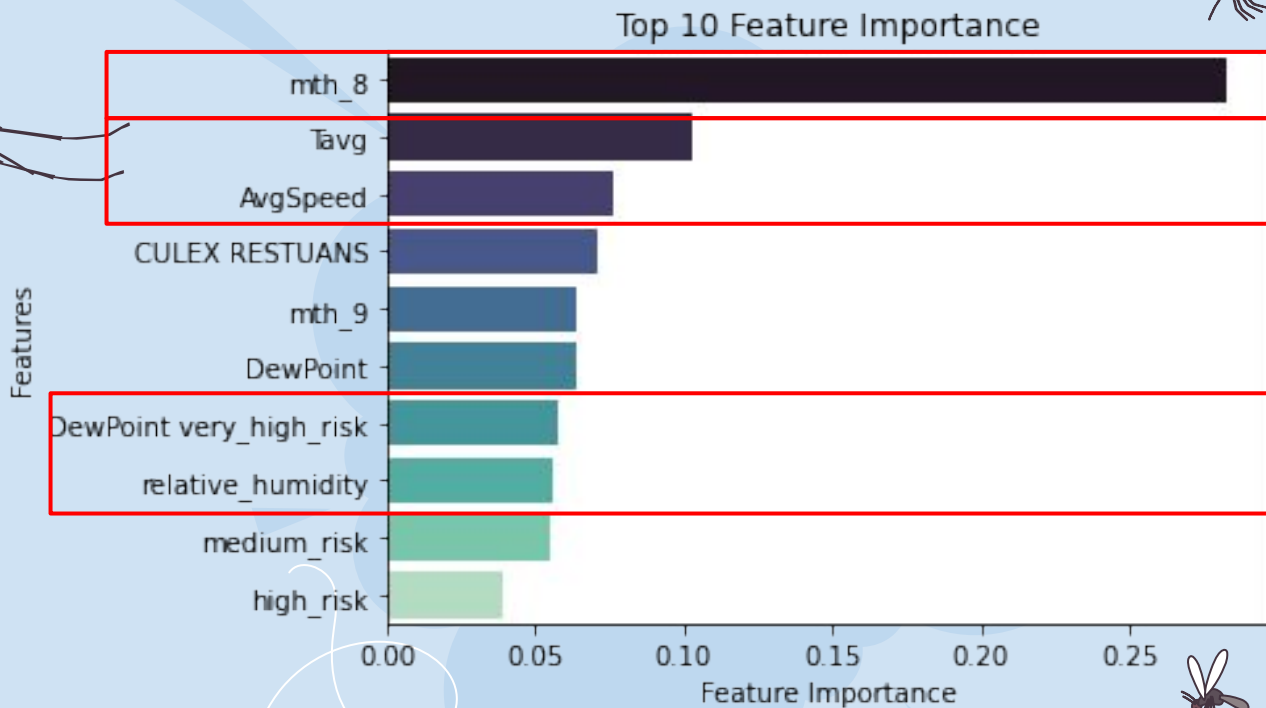


Logistic Regression	Gradient Boost	XGBoost	ADA Boost	Random Forest
81.8%	82.8%	82.9%	83.0%	81.5%
Test Score	Test Score	Test Score	Test Score	Test Score
84.5%	89.4%	89.6%	89.3%	93.7%
Train Score	Train Score	Train Score	Train Score	Train Score
<ul style="list-style-type: none"><li>• Baseline</li><li>• No parameter tuning</li></ul>	<ul style="list-style-type: none"><li>• SMOTE</li><li>• Best recall: 0.49</li><li>• Precision: 0.22</li><li>• Best F1: 0.3</li></ul>	<ul style="list-style-type: none"><li>• SMOTE</li><li>• Recall: 0.48</li><li>• Precision: 0.21</li><li>• F1: 0.29</li></ul>	<ul style="list-style-type: none"><li>• SMOTE</li><li>• Worst recall: 0.11</li><li>• Best precision: 0.25</li><li>• Worst F1: 0.15</li></ul>	<ul style="list-style-type: none"><li>• SMOTE</li><li>• Recall: 0.38</li><li>• Most overfitted model</li></ul>

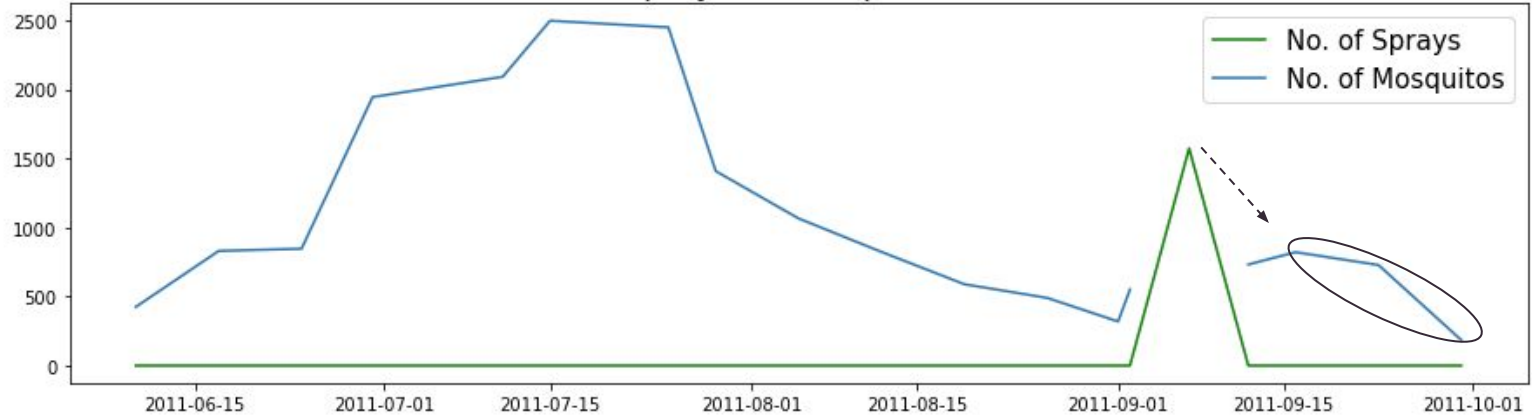
# Confusion Matrix of Actual vs Predicted WMV



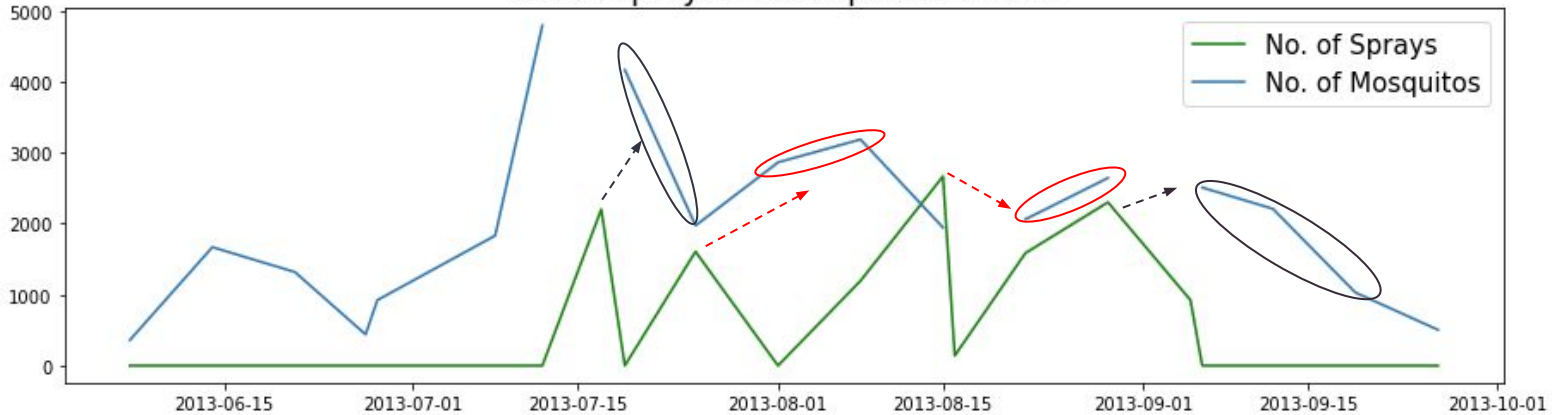
# Top 10 Feature Importance



Plot of sprays vs mosquito numbers

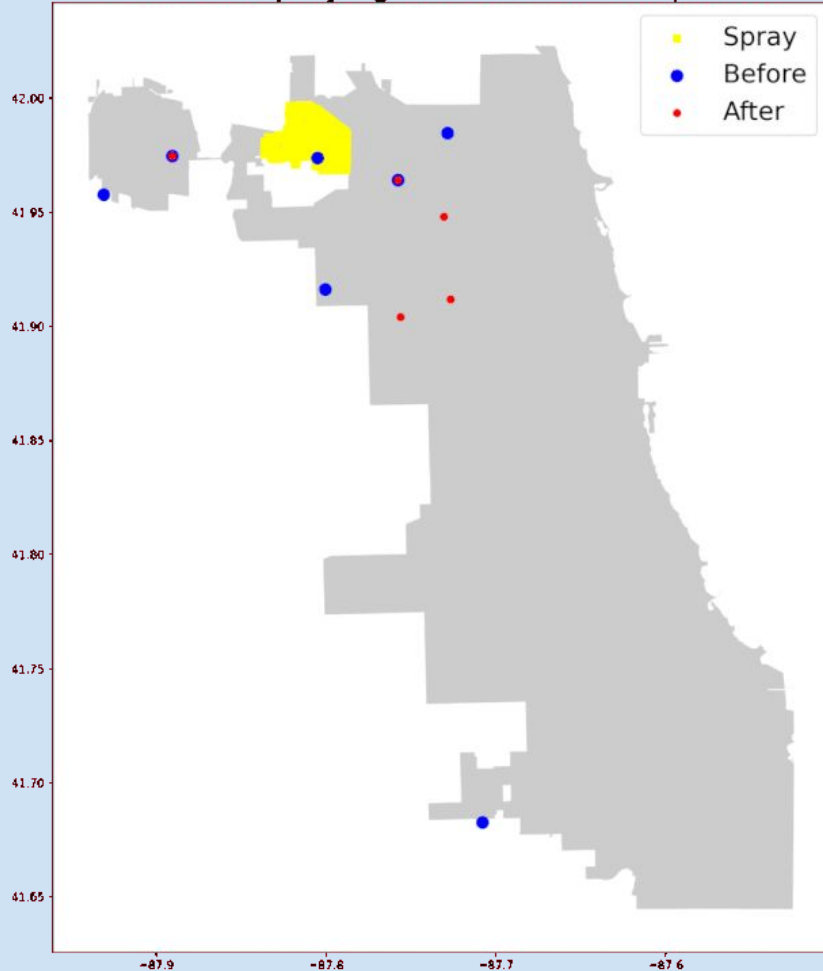


Plot of sprays vs mosquito numbers





## Effects of spraying 1 week after 7 Sep 2011



Date	Effective Areas	Ineffective Areas
7 Sep 2011	1	0
17 Jul 2013	1	1
25 Jul 2013	1	0
8 Aug 2013	2	0
15 Aug 2013	1	3
22 Aug 2013	0	3
29 Aug 2013	1	2
5 Sep 2013	0	1
<b>Total</b>	<b>7</b>	<b>10</b>

**41.18% success rate!**

# Cost Benefit Analysis

**Taking into account:**  
**Previous Sprays 41.9%**  
**< Predicted Sprays 49%**



**Spray**

**Benefit**



Reduce disease in population



Lower load on hospital services



Overall better quality of life



Approximate cost of \$4.5 million

**No Spray**

**Cost**



High Medical Expenses



Loss of Workplace Productivity



More death rates



Approximate cost of \$1 million

# Conclusion



- **Gradient Boost Classifier** has a roc-auc score of **83%** with highest recall score of **49%**
- Concentrate spray efforts during month of **August**
- Consolidate weather information on **Temperature / Wind Speed / Humidity**
- **Location** a relatively strong predictor, WNV may be **transient**



# Recommendations

<b>Adult Control</b>	Spray during (month/location/weather)
<b>Larval Control</b>	Larva-ciding water bodies
<b>Public Education</b>	Educated public on personal protection
<b>Bird Surveillance</b>	Monitor birds known to be carriers, data collected to enhance model

