

# Project 4 : West Nile Virus Prediction

General Assembly DSI 14

July 8, 2020

Song Yuan, Qi Wen, Jin, Eng Seng

# Problem Statement

The main goal here for our analytics effort is to put systems in place that reduce people's exposure to mosquitoes that carry WNV

- Gene Leynes, Data Scientist

# Questions?

How to reduce WNV  
incidences?

Where and how does  
WNV appear?

How to make prevention  
methods more effective?

Cost Benefit?

## History of West Nile Virus In Chicago



WNV Appears

West Nile Virus Exist In United States



64 Deaths

Total of 884 Cases Reported



10 & 13 Deaths



12 Deaths

**THE GOAL IS 0 DEATHS**

# TABLE OF CONTENTS

## 01 About the West Nile Virus Vectors

EDA are worth Thousand Words!

## 02 Modeling

Modeling Process and Scoring

## 03 Cost Benefit Analysis

Cost Analysis

## 04 Conclusion

Key Findings, Recommendations, Limitations & Further Research

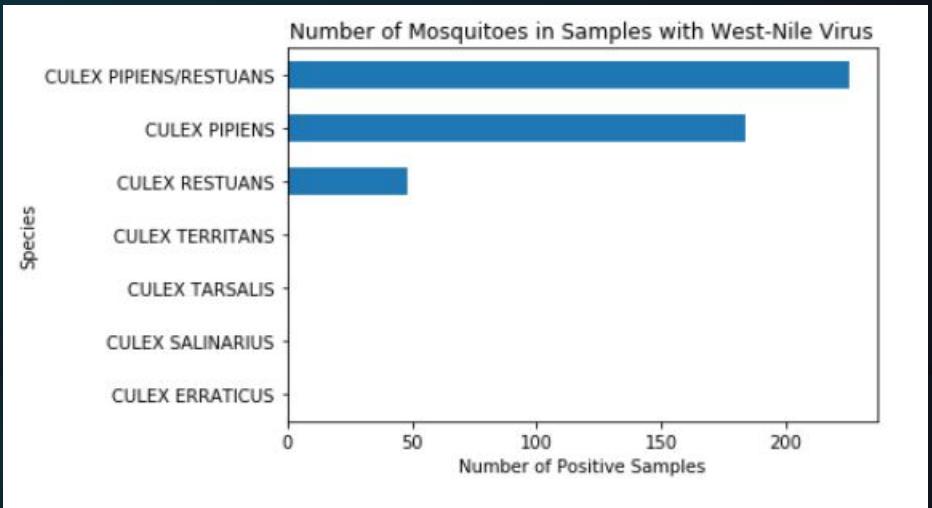
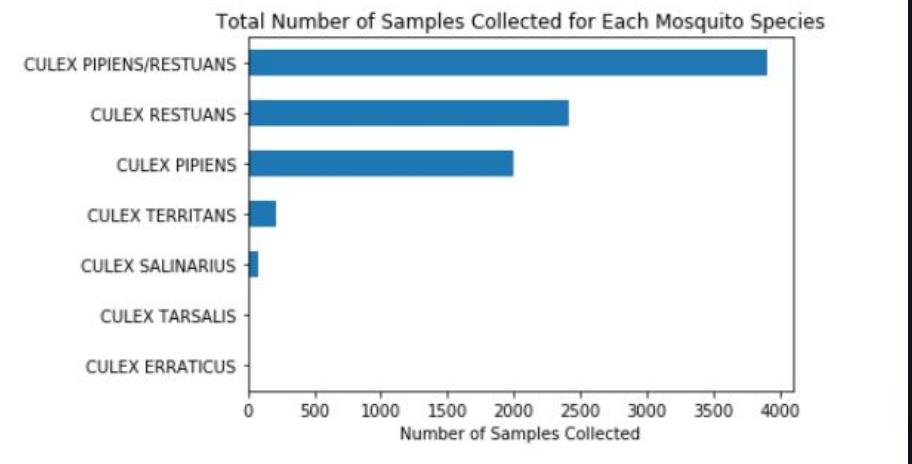
# OI

## About West Nile Virus Vectors

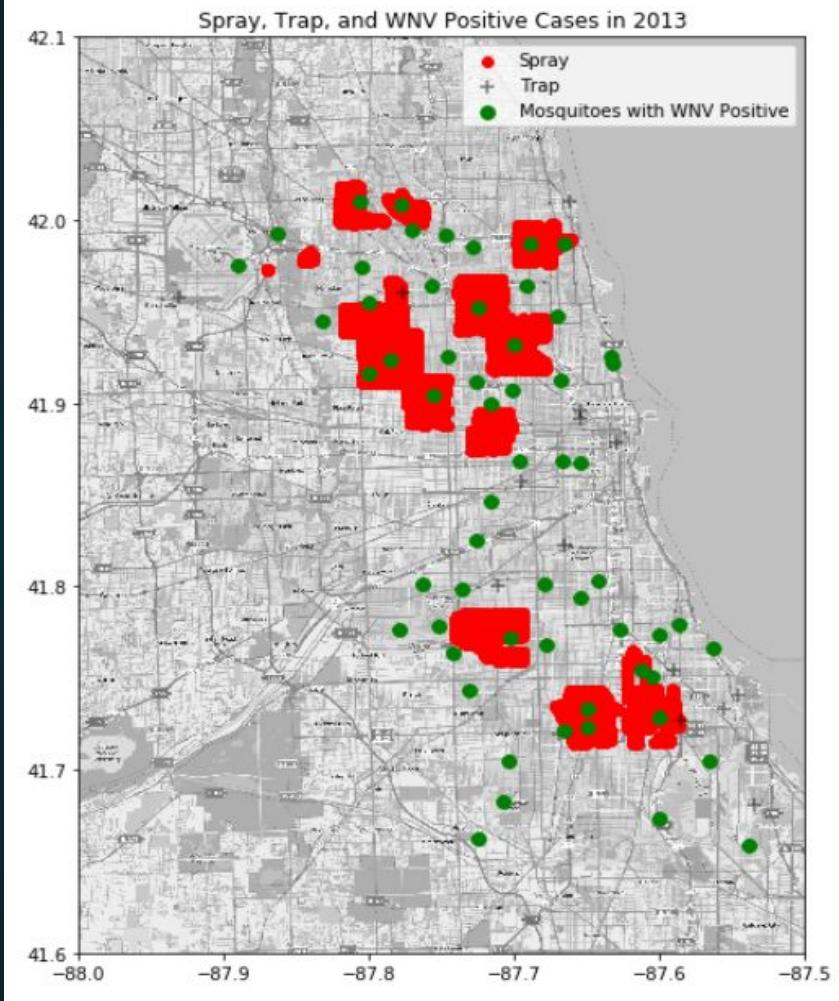
A black and white photograph of a man in a suit sitting in a car. He is looking down at a laptop computer, which is open on his lap. His right hand is on the keyboard, and his left hand is holding a white coffee cup with a lid. He is wearing a dark suit jacket, a light-colored dress shirt, and a patterned tie. A silver watch is visible on his left wrist. The background shows the interior of a car with a window and some blurred lights.

**EDA's ARE  
WORTH A  
THOUSAND  
WORDS**

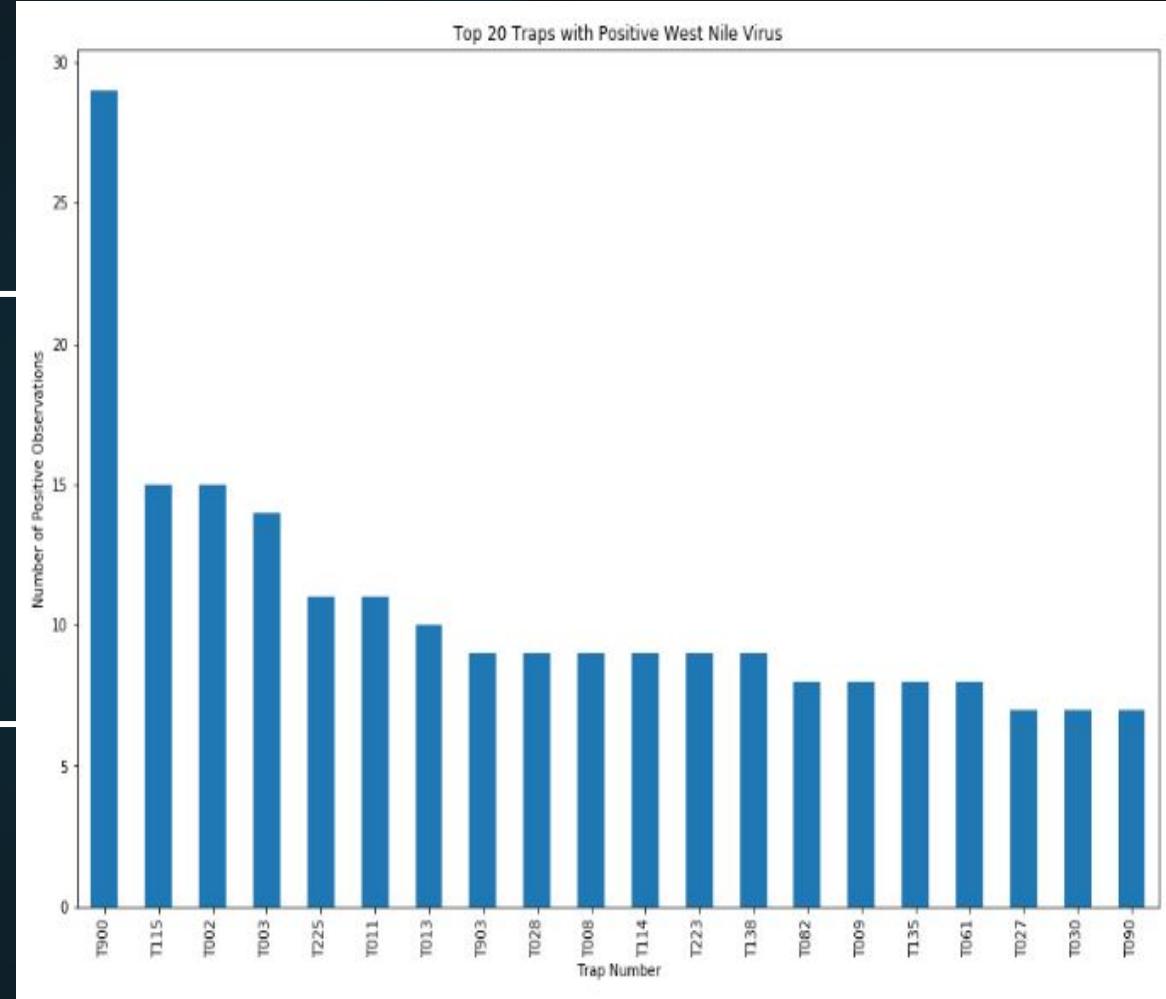
# Mosquito Species Samples Collected



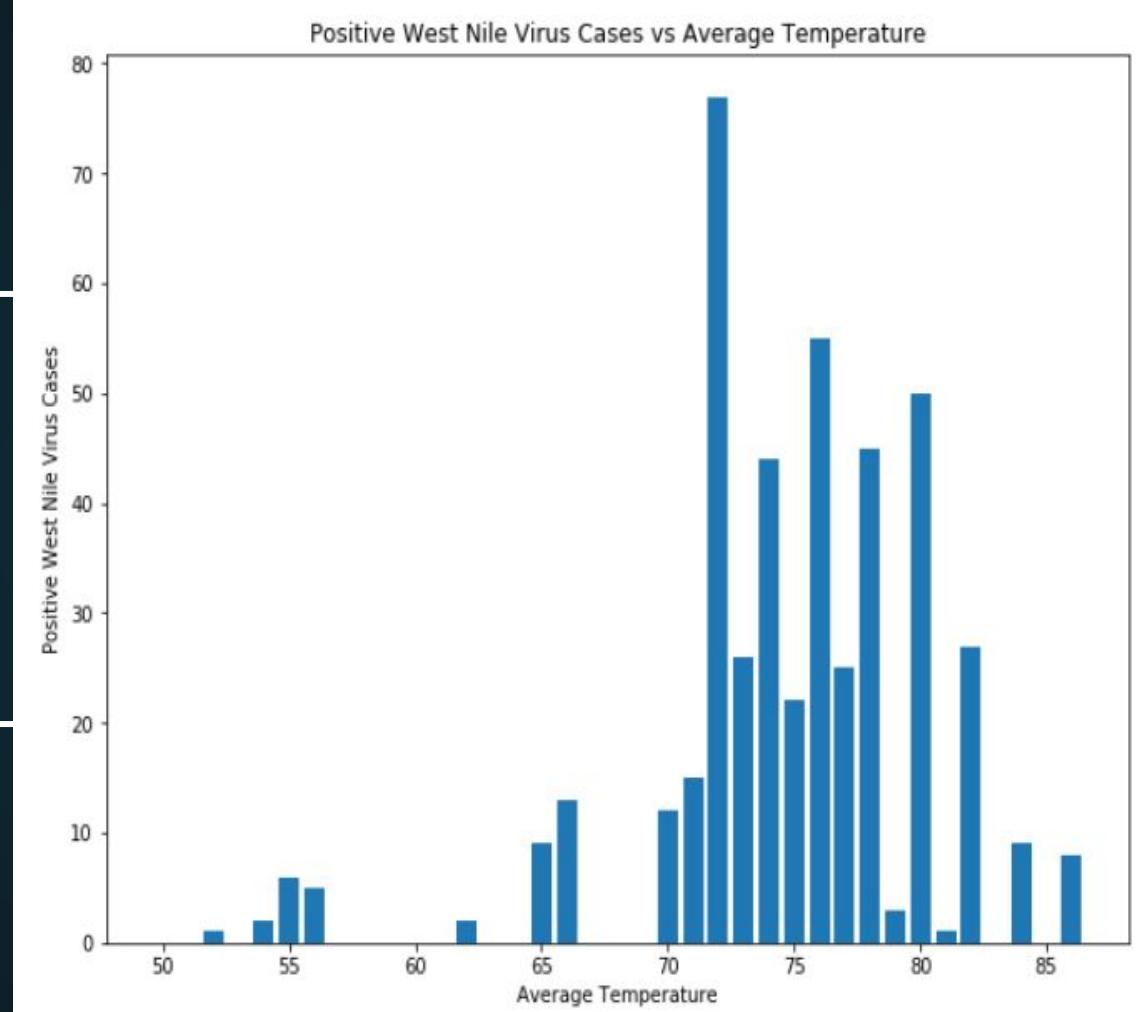
# EDA for WNV Positive Cases, Spray and Trap Mapping



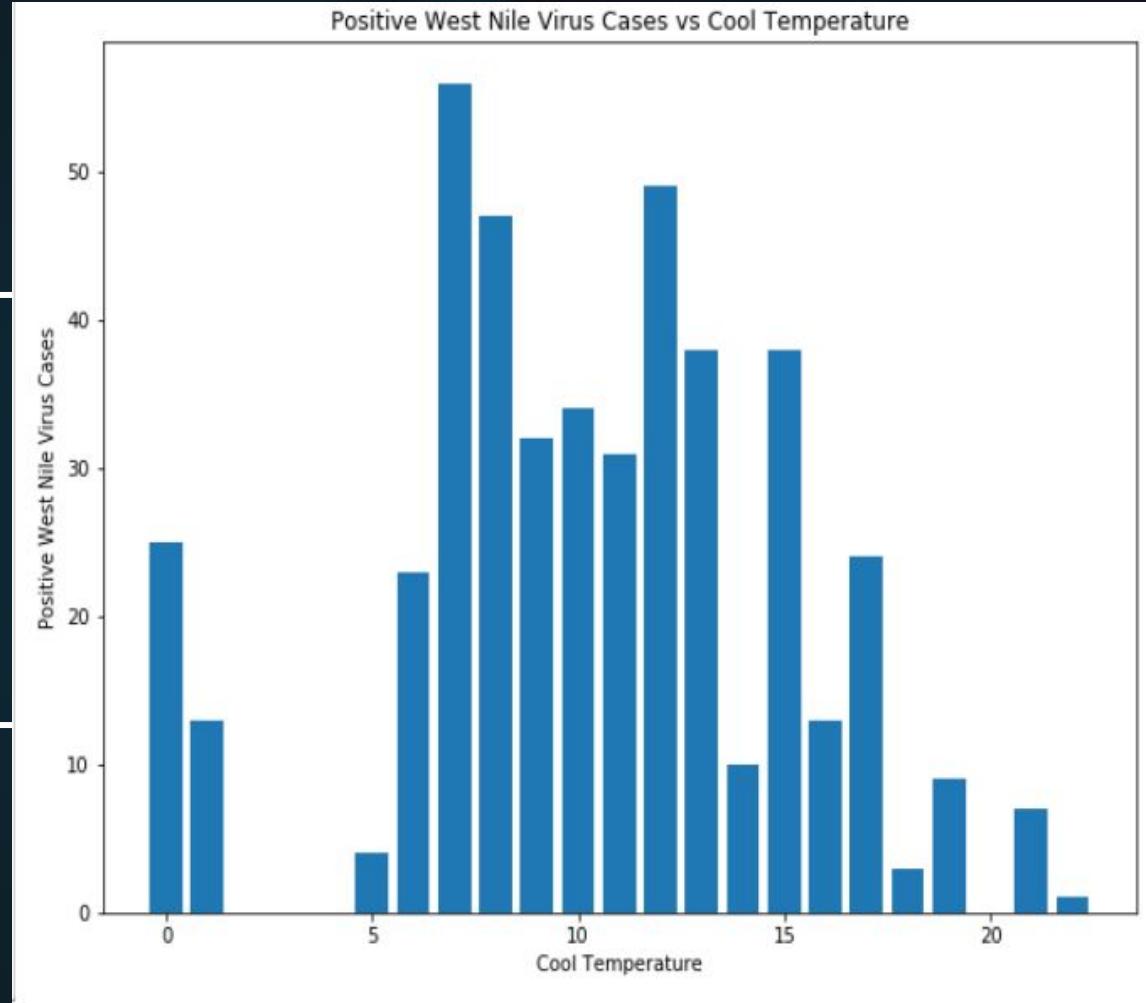
## Top 20 Traps With WNV Recorded



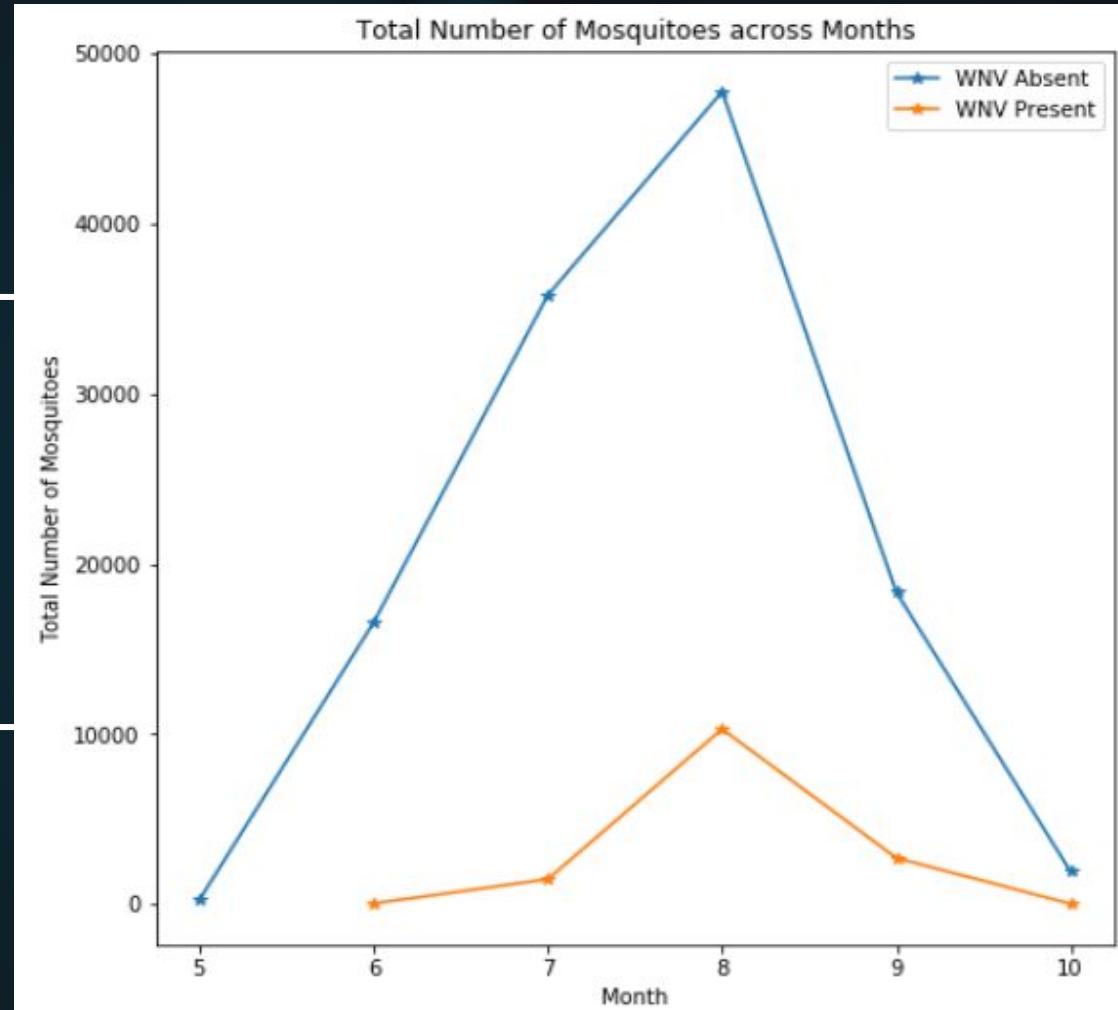
Does Temperature  
Affect the  
Existence of WNV  
??



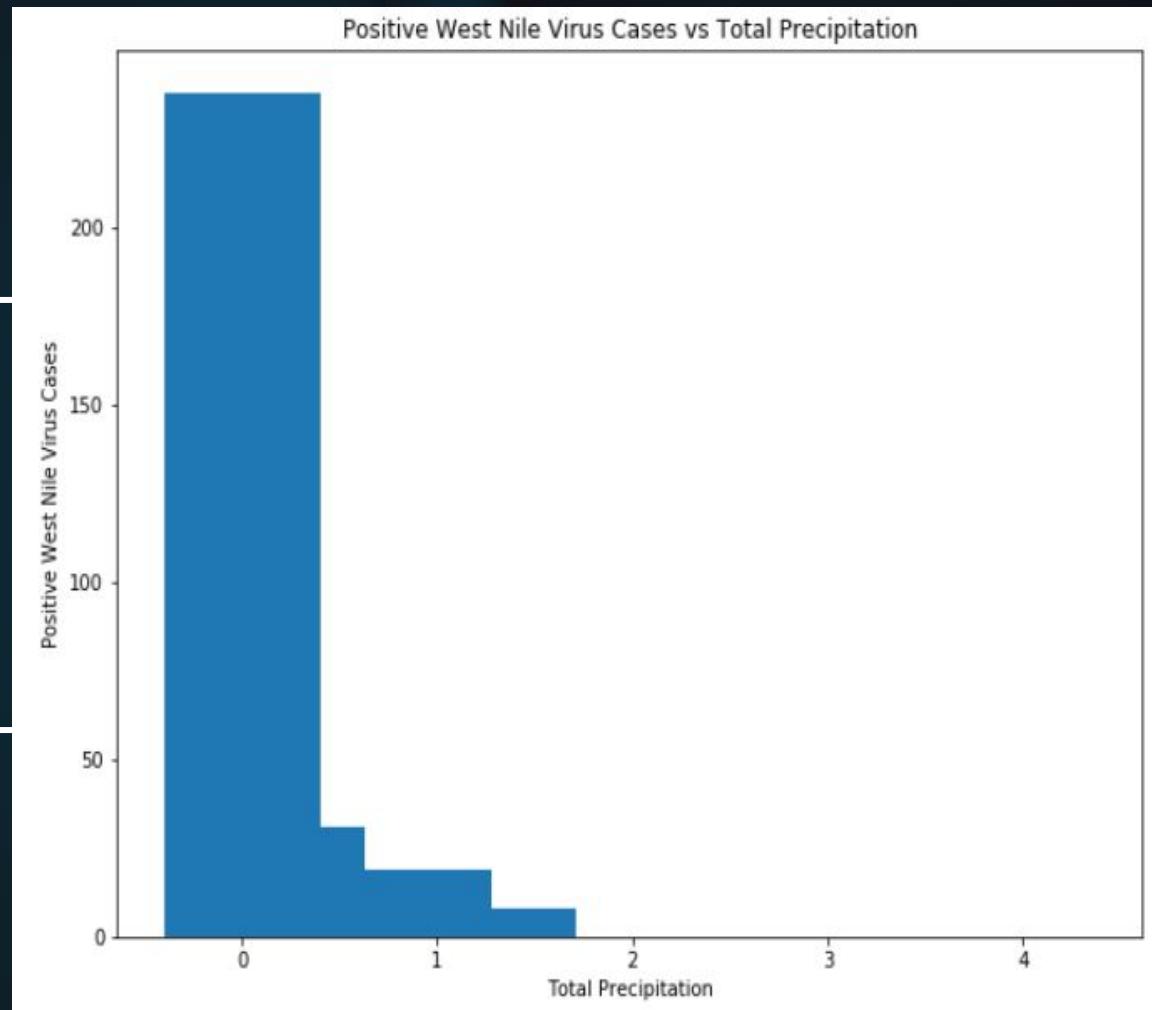
## What About Cool Temperature?



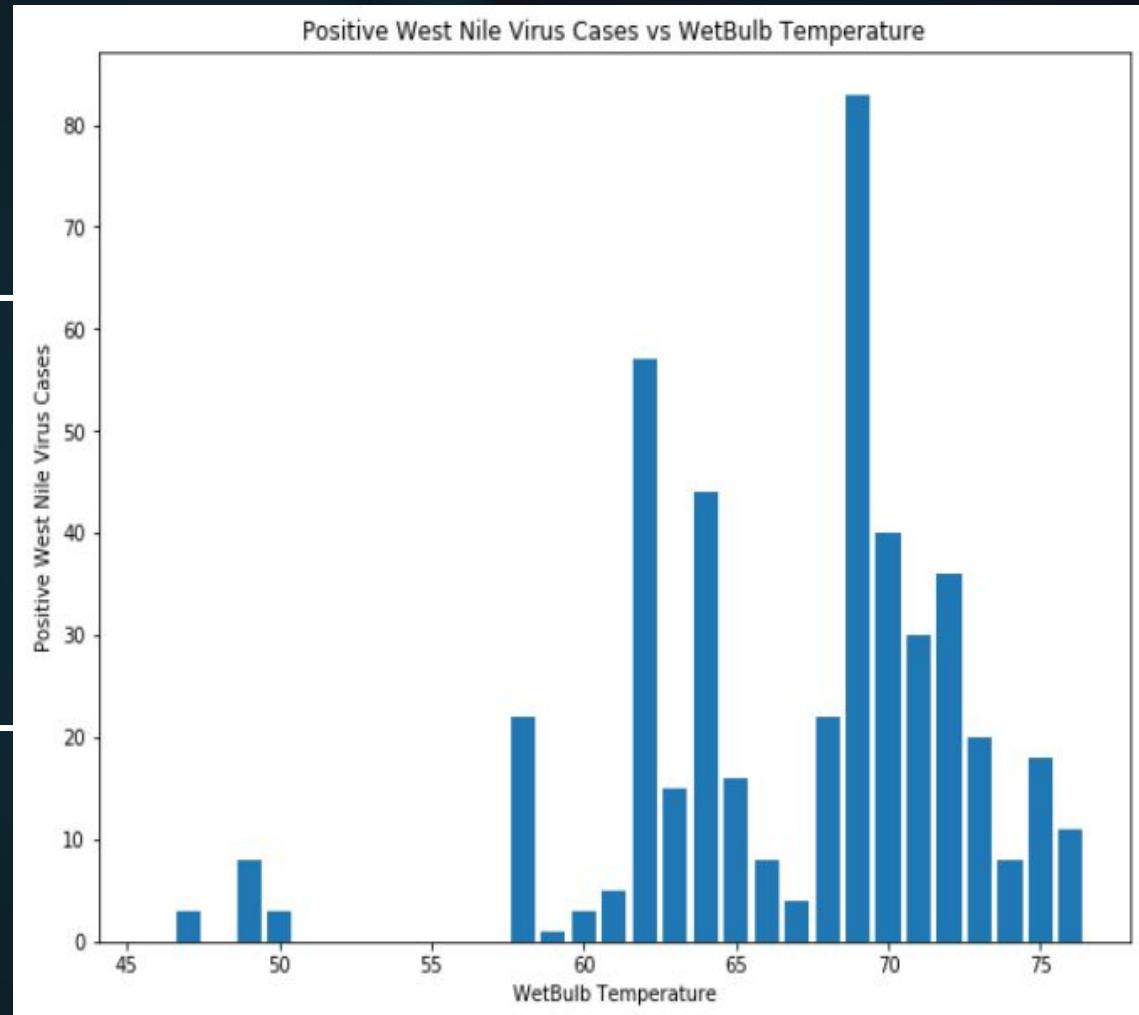
**Summer Is the  
SEASON OF  
INFECTION !!!**



**When the Weather  
precipitation is  
low , the more  
WNV IS !!!**



The higher  
humidity the more  
frequent WNV  
detected!



02

---

# Modeling



# Dataset

Initial Dataset

**Training Data** - trap locations, number of mosquitoes, presence of West Nile Virus (2007, 2009, 2011, 2013)

**Weather Data (from NOAA) -** Weather conditions from 2007 to 2014, during the months of the tests

**Spray Data** - in 2011, 2013

**Testing Data** - predict results for 2008, 2010, 2012, and 2014

Cleaning & Preprocessing

**Weather Data**

- Drop correlated columns, columns with more than 50% missing values

**Train & Test Data**

- Drop address columns
- Create new columns for total number of mosquitoes

**Spray**

- Drop column with many duplicates and missing values

Fusing Data: Match the nearest weather station to each trap

One Hot Encoding: Species and Trap number (168 features)

Final Dataset

Split train dataset into training and validation data

- Training: SMOTE oversampling to 11414 rows
- Validation: 2583 rows

Testing Set: 116293 rows

## How We Built the Models...

- The key output is score which indicating the risk that a specific point/area could test positive for WNV in particular time.
- With Variety models tested, we keep the features which will improve the predictive ability
- For each trap, the model will predict the probability the WNV would be exist

# Modeling Scoring

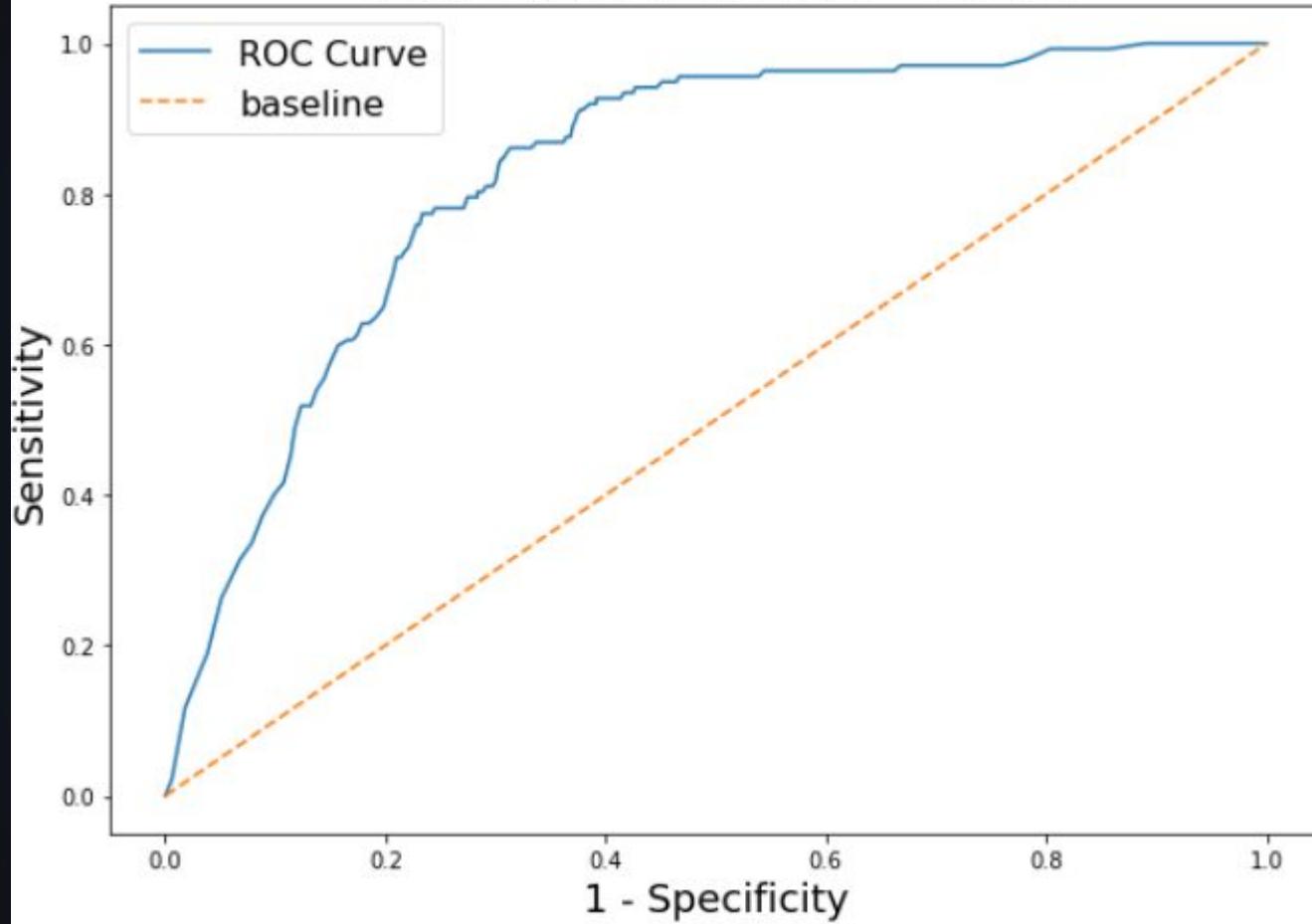
MODELS (SMOTE)	BEST PARAMS	KAGGLE PRIVATE SCORE	KAGGLE PUBLIC SCORE
Logistic Regression (lbfgs)	C = 1, Penalty = l2, max_iter = 1000	0.5719	0.5022
Logistic Regression (liblinear)	C = 1, Penalty = l1, max_iter = 1000	0.5719	0.5022
Decision Tree	max_depth = 10, min_samples = 4, min_samples_split = 20	0.5505	0.5873
Random Forest	max_depth = 20, max_leaf_nodes = 20, min_samples_leaf = 10, n_estimators = 100	0.5765	0.5882
AdaBoost	learning_rate = 1.0, n_estimators = 100	0.5463	0.5481
XGBoost	eval_metric = 'auc', scale_pos_weight = 18, subsample = 0.5, eta = 0.2	0.68208	0.6974

# Results

The most predictive feature

Features	Importance
Month	<b>0.036464</b>
Trap_T095	<b>0.028762</b>
Trap_T225	<b>0.018460</b>
CULEX TERRITANS	<b>0.017817</b>
DATE	<b>0.017134</b>
Trap_T094	<b>0.016628</b>
Trap_T158	<b>0.015939</b>
Trap_T080	<b>0.015407</b>
Trap_T063	<b>0.015346</b>
Trap_T102	<b>0.015210</b>

### ROC Curve with AUC = 0.828



 **ROC CURVE  
WITH AUC**

**Scoring :  
0.828**

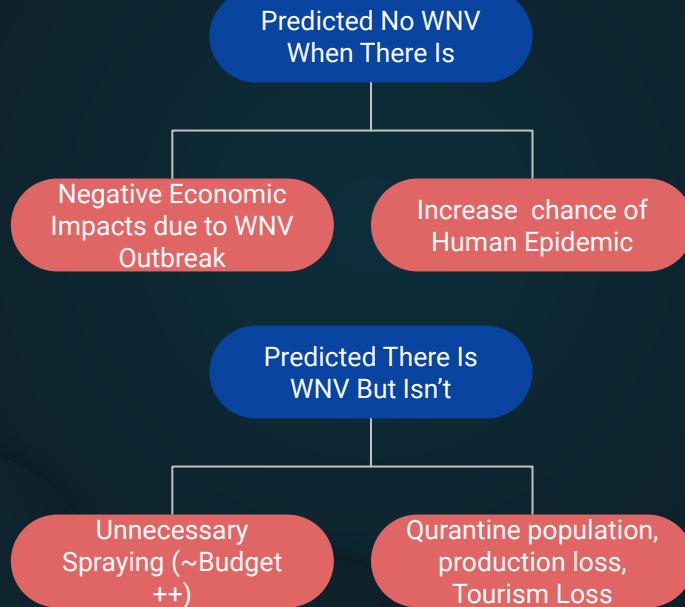
03

## Cost Benefit Analysis

What are the real cost and benefits  
of vector control?



## The Confusion Matrix



# **906 Cases**

Reported From Chicago from Year 2005 to  
2016

# **\$ 12,977,000**

Estimated Medical Cost from Year 2005 to  
2016

# **\$ 702,000**

49 Cases in 2016. Costing \$ 702k in medical expenses  
and lost productivity



# **148,510 Acres**

Chicago city Area

# **\$ 137,000 estimated**

Estimated for worst case scenario of  
spraying whole city

# **5 Times**

Medical Expenses & Loss Productivity cost incurred 5  
times more expensive than spraying in the city !

04

---

## Conclusion

## Key Findings

- Mosquito numbers follow the natural onset of the season, begin in May and end in October.
- Peak at August and decreasing at October (died / hibernation due to winter)
- Spray Operations by City of Chicago and Department of Public Health do not have noticeable effect.
- Past spray operations had missed locations that reported the most number of positive cases
- Strong Correlation between temperature and precipitation with number of WNV positive!
- Most WNV positive samples found between 71F - 80F, zero precipitation.

## **Recommendations**

- Better spraying regime that is informed by weather and trap data.
- Concentrate efforts to combat WNV prone areas which tend to fall into two categories outline. Inner suburbs fairly high income, 1940s - 1960s housing and Low income, young in City.
- More public awareness campaigns. Urging people to minimize exposure with necessary precautions steps.

- The dataset was extremely imbalanced and we had to generate synthetic samples to balance the dataset
- Spraying data was limited to only 2 years with no indications of dosage
- Due to computational power limit, we had only tuned limited combinations of hyperparameters
- Other mosquito control efforts have not been accounted for in our model determine the actual effectiveness of spraying
- For a more complete analysis on economical and costings, we should collect the data in terms of healthcare expenses and mosquito control efforts.

## **Limitations & Further Research**

# THANKS!

ANY QUESTIONS ???