# Project 4 -West Nile Virus

Alex, Elliot, Mak & Robby

## Table of Contents

1. INTRO & EDA

Mak

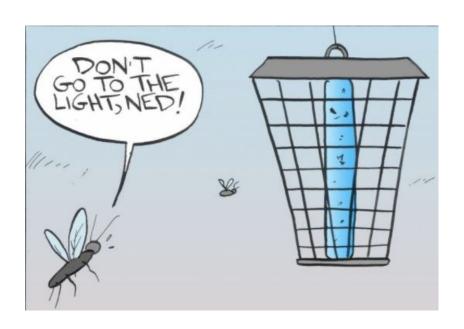
PREPROCESSING,

- 2. **FEATURE SELECTION & MODELING**Alex
- 3. **MODEL EVALUATION** Robby
- 4. **COST BENEFIT ANALYSIS** Elliot

# An Overview of the Problem



## Background



- In reaction to an earlier 2002 outbreak, during 2004 the Chicago Department of Public Health(CDPH) set up a surveillance and control system to trap and test mosquitoes for the presence of WNV.
- Public health workers in Chicago setup mosquito traps scattered across the city.
   These traps collect mosquitos, and the mosquitos are tested for the presence of West Nile virus.

## Problem Statement

As data scientists from a data consultancy firm engaged by the Chicago Department of Public Health, our goal is to build a fully predictive model to be used for proactive mosquito management.

#### The effects are two-fold:

- Saving costs in terms of vector management and spraying specific traps predicted by our model which will be more effective and targeted. This helps the department to stay within their budget.
- 2. The proactive and predictive approach will prevent wnv infections from reaching outbreak levels. The benefits accrued for this is the cost savings from human and medical costs that might otherwise be incurred from wnv infections.

## Data Exploration & EDA

## Data Sources

```
weather = pd.read_csv('../datasets/weather.csv')
weather.head()
```

	Station	Date	Tmax	Tmin	Tavg	Depart	DewPoint	WetBulb	Heat	Cool	 CodeSum	Depth	Water1	SnowFall	PrecipTotal	StnPressure	SeaLevel	R
0	1	2007- 05-01	83	50	67	14	51	56	0	2		0	М	0.0	0.00	29.10	29.82	
1	2	2007- 05-01	84	52	68	М	51	57	0	3		М	М	М	0.00	29.18	29.82	
2	1	2007- 05-02	59	42	51	-3	42	47	14	0	 BR	0	М	0.0	0.00	29.38	30.09	
3	2	2007- 05-02	60	43	52	М	42	47	13	0	 BR HZ	М	М	М	0.00	29.44	30.08	
4	1	2007- 05-03	66	46	56	2	40	48	9	0		0	М	0.0	0.00	29.39	30.12	

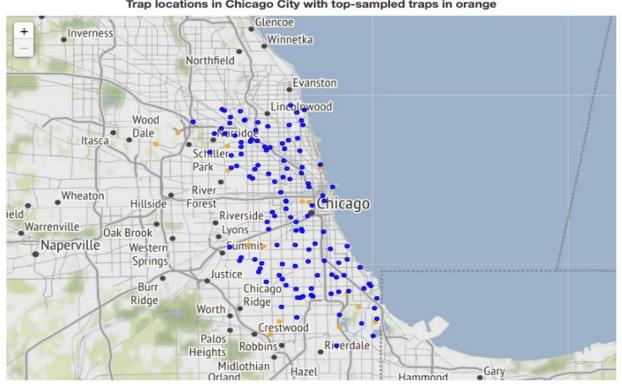
spray = pd.read\_csv('../datasets/spray.csv')
spray.head()

	Date	Time	Latitude	Longitude
0	2011-08-29	6:56:58 PM	42.391623	-88.089163
1	2011-08-29	6:57:08 PM	42.391348	-88.089163
2	2011-08-29	6:57:18 PM	42.391022	-88.089157
3	2011-08-29	6:57:28 PM	42.390637	-88.089158
4	2011-08-29	6:57:38 PM	42.390410	-88.088858

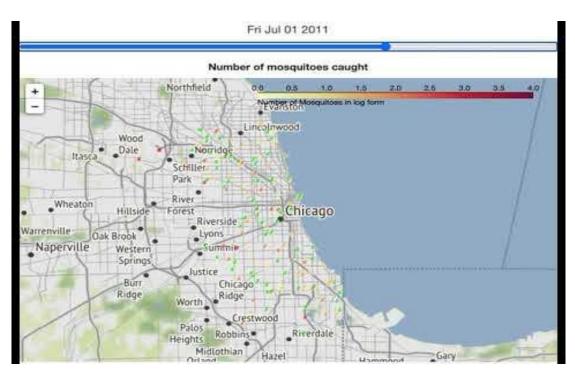
<pre>train = pd.read_csv('/datasets/train.csv') train.head()</pre>						test = pd.read csv('/datasets/test.csv')											
							head()	uu_cov( uucubecb/ cc	,								
1100 North Coak CULEX Park CULEX Park CULEX Chicago, Chicago, Li 6034,	N OAK RK AVE T002 4100 N OAK PARK Chica	AVE. 41.954690 -87.800991	9	1	c	ld	Date	Address	Species		100000000000000000000000000000000000000	112000000	4400 N OAK BARK N/E			AddressAccuracy	
4100 North Oak						0 1	2008- 06-11	4100 North Oak Park Avenue, Chicago, IL 60634,	CULEX PIPIENS/RESTUANS	41	N OAK PARK AVE	T002	4100 N OAK PARK AVE, Chicago, IL	41.95469	-87.800991	9	
1 2007- Park CULEX RESTUANS 41 P. Chicago, IL 60634,	N OAK T002 4100 N OAK PARK RK AVE Chica	AVE, 41.954690 -87.800991 go, IL	9	1	Ç	1 2	2008- 06-11	4100 North Oak Park Avenue, Chicago, IL 60634,	CULEX RESTUANS	41	N OAK PARK AVE	T002	4100 N OAK PARK AVE, Chicago, IL	41.95469	-87.800991	9	
6200 North Mandell 2 2007 - Avenue, CULEX RESTUANS 62 M	N ANDELL T007 6200 N MANDELL	AVE, 41.994991 -87.769279	9	1	c	2 3	2008- 06-11	4100 North Oak Park Avenue, Chicago, IL 60634,	CULEX PIPIENS	41	N OAK PARK AVE	T002	4100 N OAK PARK AVE, Chicago, IL				
IL 60646, USA	AVE Cnica	go, IL				3 4	2008- 06-11	4100 North Oak Park Avenue, Chicago, IL 60634,	CULEX SALINARIUS	41	N OAK PARK AVE	T002	4100 N OAK PARK AVE, Chicago, IL	41.95469	-87.800991	9	
7900 West Foster 3 2007- Avenue, CULEX 79 05-29 Chicago, PIPIENS/RESTUANS 79	W TO15 7900 W FOSTER AVE Chica	AVE. 41.974089 -87.824812 go, IL	8	1	c	4 5	2008- 06-11	4100 North Oak Park Avenue, Chicago, IL 60634,	CULEX TERRITANS	41	N OAK PARK AVE	T002	4100 N OAK PARK AVE, Chicago, IL	41.95469	-87.800991	9	

## Traps Analysis

Trap locations in Chicago City with top-sampled traps in orange



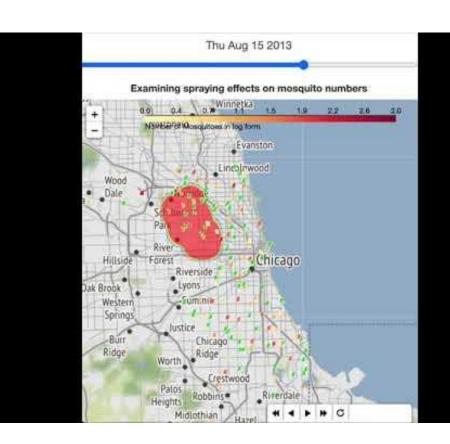
## Traps Analysis



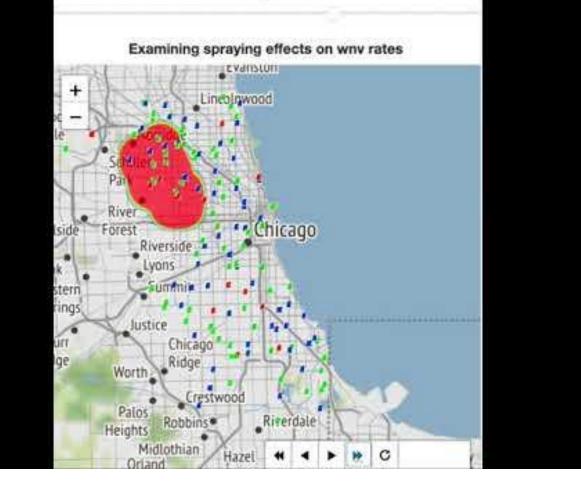
## Spray Effectiveness



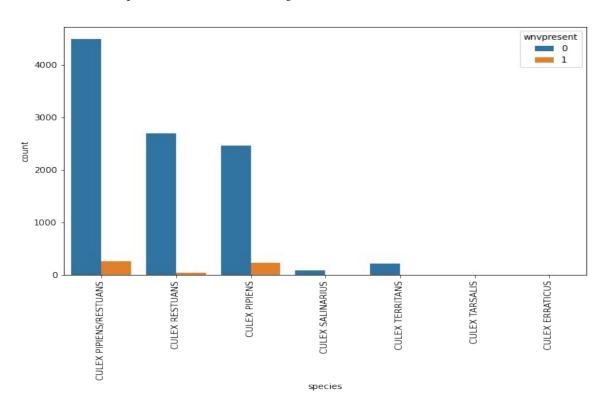
Increase (Num Mosquitoes): 9
Decrease (Num Mosquitoes): 18
No change (Num Mosquitoes): 1



#### Thu Aug 15 2013



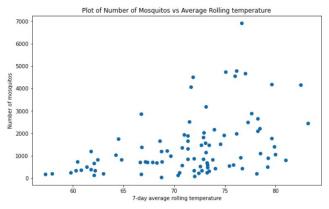
## Mosquito Analysis

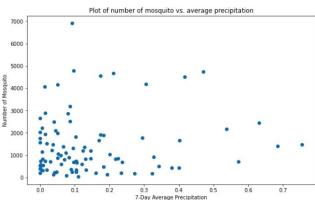


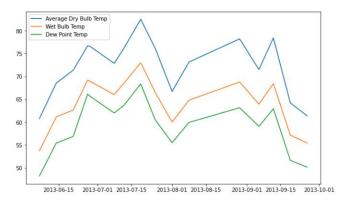
#### Main carriers for WNV:

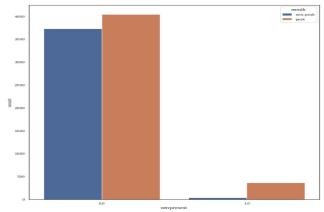
- Culex Pipiens / Restuans
- Culex Restuans
- Culex Pipiens

## Weather Analysis









- L. INTRO & EDA Mak
- PREPROCESSING,
- 2. **FEATURE SELECTION & MODELING**Alex
- 3. MODEL EVALUATION Robby
- 4. **COST BENEFIT ANALYSIS** Elliot



# Preprocessing

## Weather Preprocessing

- 2 Weather Stations
  - Took the average
  - Took the maximum temperature
  - Took the minimum temperature
- Feature Engineering -
  - DayLength
  - Weather Rolling 7
  - Weather Rolling 14

#### Handling of the DewPoint, WetBulb, StnPressure & AvgSpeed results

DewPoint: Take the average result of the 2 stations.

WetBulb: Take the average result of the 2 stations.

StnPressure: Take the average result of the 2 stations.

AvgSpeed: Take the average result of the 2 stations.

#### Handling of the Tmax, Tmin, Tavg results

Tmax: Take the maximum temperature of the 2 stations.

Tmin: Take the minimum temperature of the 2 stations.

Tavg: Take the average temperature of the 2 stations.

#### DayLength: Create a new feature using Sunset - Sunrise.

```
weather_rolling_7 = weather_final.rolling(7).mean()

weather_rolling_7.columns = [col_name + '_roll_7' for col_name in weather_final.columns]

weather_rolling_14 = weather_final.rolling(14).mean()

weather_rolling_14.columns = [col_name + '_roll_14' for col_name in weather_final.columns]
```

## Preprocessing

#### **Traps**

- Several numbers in the test file were not present in the training file.
- These traps were replaced with known traps nearby

#### **Additional Feature Engineering**

- One hot encoded Species, Trap & Month
- Created feature wnv\_species
- Created feature wnv\_traps (Binary classification for traps with WNV)

```
Nearest trap for T090A is T086 at 1658.96 meters away
Nearest trap for T090B is T090 at 438.74 meters away
Nearest trap for T090C is T086 at 2508.67 meters away
Nearest trap for T200A is T099 at 608.96 meters away
Nearest trap for T128A is T099 at 1152.9 meters away
Nearest trap for T200B is T221 at 1459.9 meters away
Nearest trap for T218A is T001 at 281.18 meters away
Nearest trap for T218C is T222 at 1147.97 meters away
Nearest trap for T218B is T001 at 857.3 meters away
Nearest trap for T002A is T017 at 678.5 meters away
Nearest trap for T002B is T002 at 1021.28 meters away
Nearest trap for T234 is T005 at 1212.58 meters away
Nearest trap for T065A is T065 at 278.81 meters away
```

# Modeling

## Modeling

- With positive WNV ( or 1) occurrence as the minority class(approx.5%) typical accuracy metrics would be undesirable.
- AUC-ROC scoring was used as our evaluation for scoring
- Additionally, due to the potentially high costs of false negatives in this case, we also focused on the Recall Score.

#### **Baseline Accuracy**

```
y.value_counts(normalize=True)
```

0.0 0.949927 1.0 0.050073

Name: wnvpresent, dtype: float64

## Modeling

- During earlier GridSearch(es), we found that boosting models gave better scores (approx. AUC =~0.86)
- Thus we decided to concentrate our efforts on boosting models

Model Prep using GridSearch with scaling and SMOTE (Gradient Boost, Adaboost and Extreme Gradient Boost)

```
pipe 1 = Pipeline([
    ('scale', StandardScaler()),
    ('sampling', SMOTE(random state=42)),
    ('gboost', GradientBoostingClassifier(random state=42))
1)
pipe 2 = Pipeline([
    ('scale', StandardScaler()),
    ('sampling', SMOTE(random state=42)),
    ('ada', AdaBoostClassifier(random state=42))
1)
pipe 3 = Pipeline([
    ('scale', StandardScaler()),
    ('sampling', SMOTE(random state=42)),
    ('xgb', XGBClassifier(random state=42))
1)
```

## Modeling (Gradient Boost & Adaboost)

2457

0.86

```
Best GridSearchCV AUC: 0.839
Training AUC on best params: 0.871
Validation AUC on best params: 0.861
```

Scoring Report for: Gradient Boosting

0.94

weighted avg

precision recall f1-score support 0.0 0.98 0.82 0.89 2334 0.18 0.72 0.28 123 1.0 0.82 2457 accuracy 0.77 0.59 2457 macro avq 0.58

0.82

Best GridSearchCV AUC: 0.834
Training AUC on best params: 0.8

Training AUC on best params: 0.869 Validation AUC on best params: 0.872

Confusion Matrix for: Ada Boosting [[2020 314] [ 43 80]]

Scoring Report for: Ada Boosting precision recall f1-score support 0.0 0.98 0.87 0.92 2334 1.0 0.20 0.65 0.31 123 0.85 2457 accuracy 0.59 0.76 0.61 2457 macro avq weighted avg 0.94 0.85 0.89 2457

# Modeling(xGB)

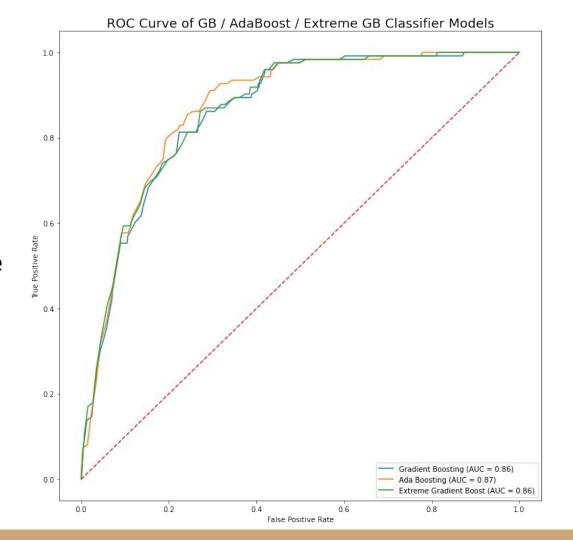
Best GridSearchCV AUC: 0.838 Training AUC on best params: 0.872 Validation AUC on best params: 0.864

Confusion Matrix for: Extreme Gradient Boost [[1827 507] [ 29 94]]

Scoring Report for: Extreme Gradient Boost precision recall f1-score support 0.0 0.98 0.78 0.87 2334 1.0 0.16 0.76 0.26 123 0.78 2457 accuracy 0.77 0.57 2457 macro avg 0.57 weighted avg 0.94 0.78 0.84 2457

## Model Evaluation

- All 3 models had **very similar ROC-AUC** (0.86-0.87)
- To understand which model to choose, we need to evaluate the misclassifications in our prediction data.

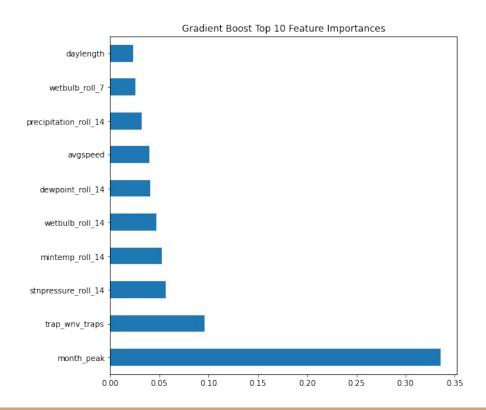


- L. INTRO & EDA Mak
- PREPROCESSING,
- 2. **FEATURE SELECTION & MODELING**Alex
- 3. **MODEL EVALUATION** Robby
- 4. **COST BENEFIT ANALYSIS** Elliot



# Model Evaluation & Misclassifications Analysis

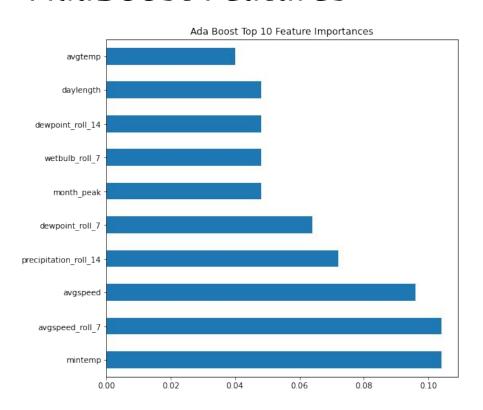
## Gradient Boost Features



#### **Observations**

- peak month was by far the most important feature that the model split on
- WNV\_traps was ranked 2nd but mosquito species was not in the top 10 features
- 14-day rolling average weather features dominate the rest of the top 10

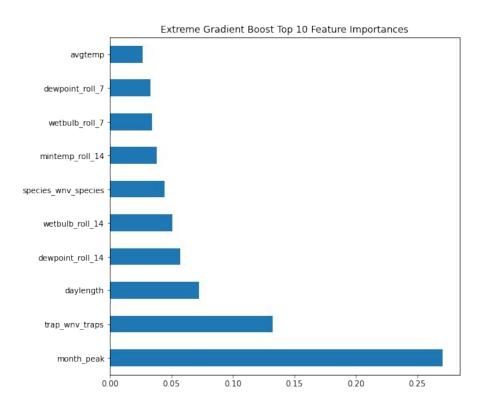
### AdaBoost Features



#### **Observations**

- Minimum temperature & 7-day average (wind) speed were the 2 most important features that the model split on
- **Combination of weather features** dominate the rest of the top 10
- Peak months was ranked 6th, the only binary feature in the top 10

## XGB Features



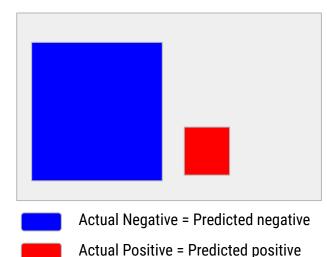
#### **Observations**

- Similar to GB model, peak months and wnv traps were the 2 most important features
- Only model to include **Mosquito species** (ranked 6th) in its top 10 features
- 7-day and 14-day rolling average weather features dominate the remaining features

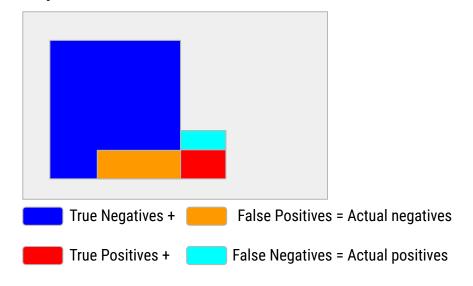
## Misclassification Analysis

Let's use set theory (venn diagrams) to imagine what a perfect model will be like:

#### **Perfect Model**



#### **Imperfect Model**



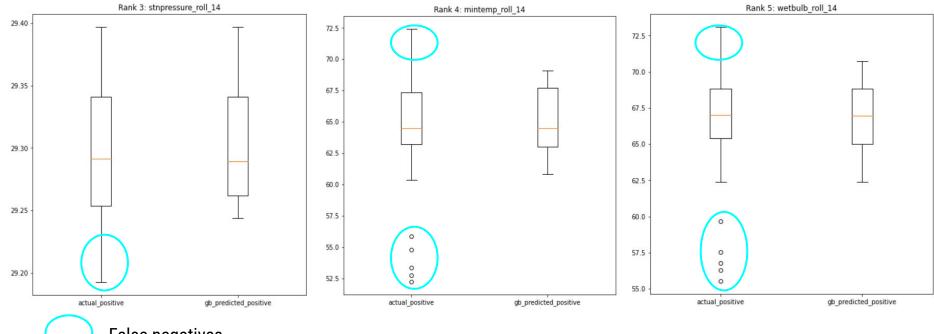
## Misclassification Analysis

Hence we would like to compare the dataset characteristics between:



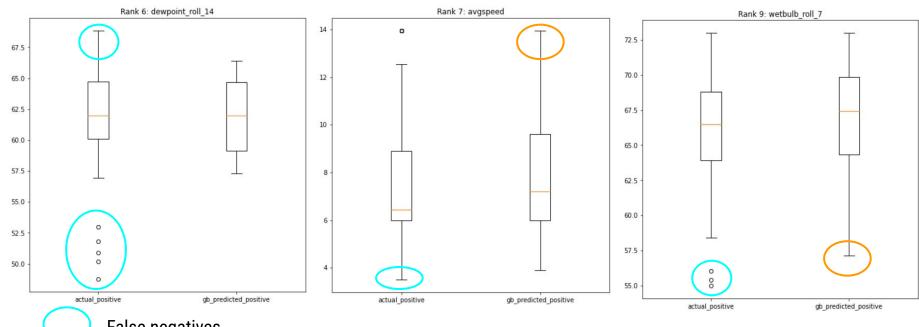
- If our tree-based models were perfect, the splits would have perfectly separated the actual positives and negatives
- Due to the imperfection, false positives (orange box) were included while false negatives were excluded (empty baby blue box) from the predicted positives

## Misclassification Analysis: GB Model



False negatives

## Misclassification Analysis: GB Model

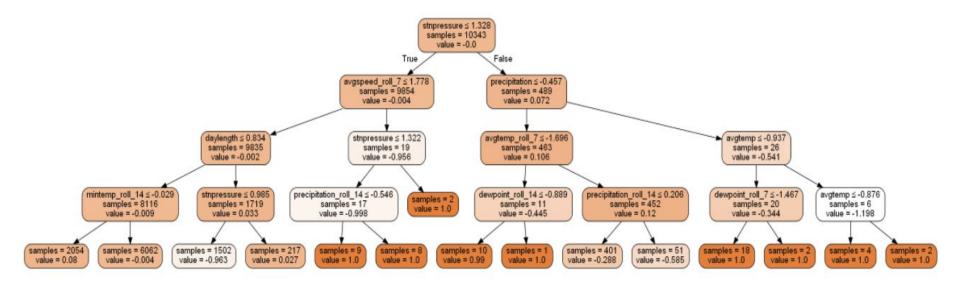


False negatives

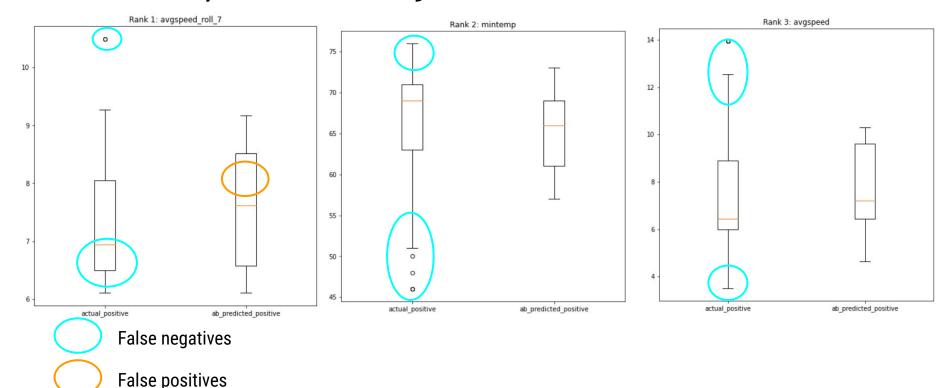
False positives

## Misclassification Analysis: GB Model

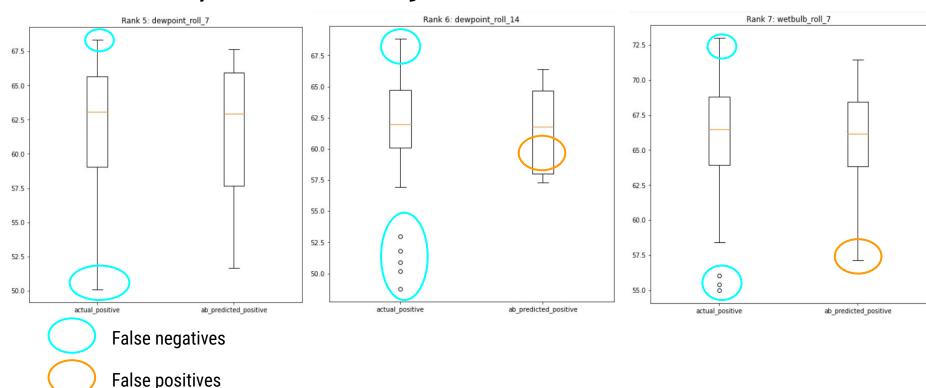
Sample tree example from GB model:



## Misclassification Analysis: AdaBoost Model

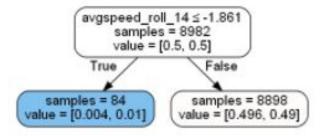


## Misclassification Analysis: AdaBoost Model



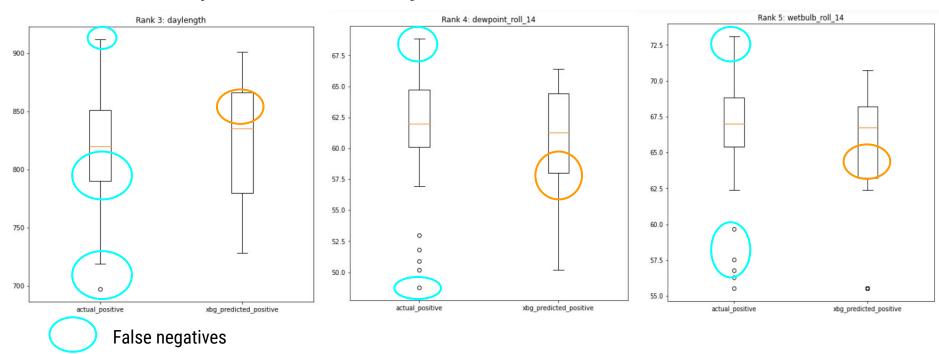
## Misclassification Analysis: AdaBoost Model

Sample tree example from AdaBoost model:

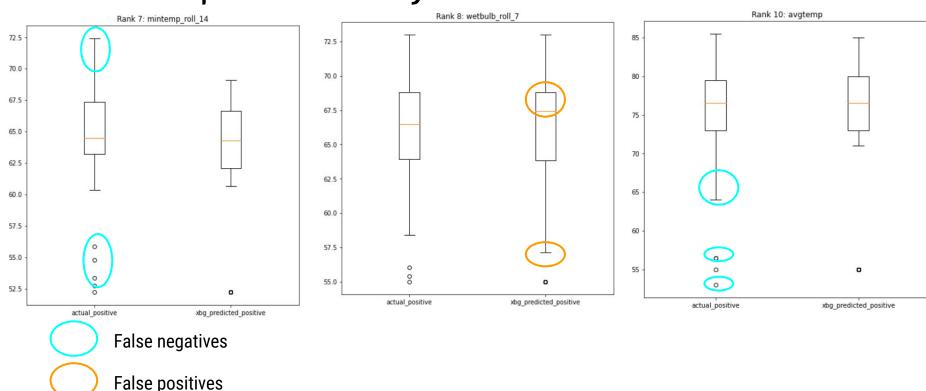


# Misclassification Analysis: XGB Model

False positives

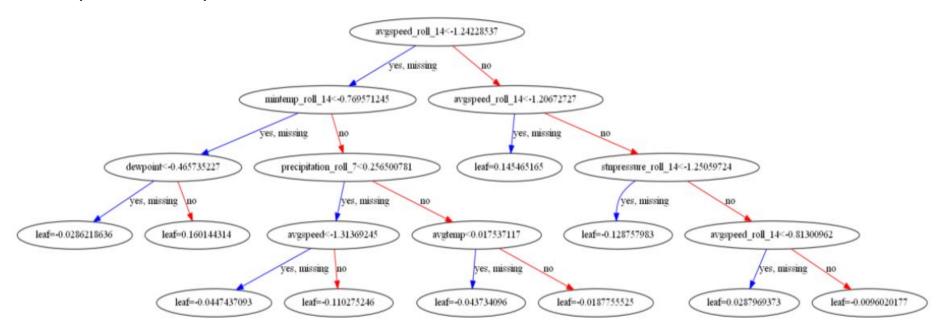


# Misclassification Analysis: XGB Model



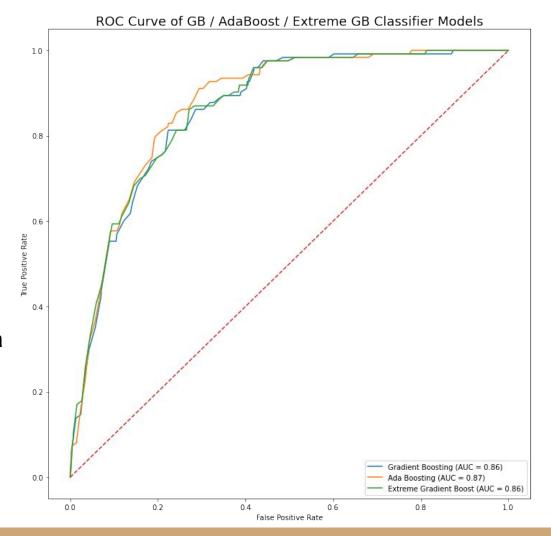
# Misclassification Analysis: XGB Model

Sample tree example from XGB model:



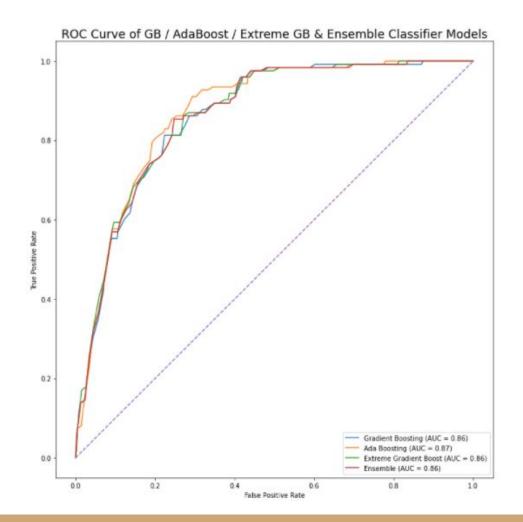
### Model Evaluation

- All 3 models had very similar ROC-AUC (0.86-0.87)
- Variety in each model's top 10 features
- Decided to have an ensemble of all 3 boosting models and build a voting classifier model



# Production model - Voting Classifier

- Voting Classifier ROC-AUC = 0.86
- Very similar scores to our 3 boosting models
- The voting classifier averages out the shortfalls of each model



# Models comparison

	TN	FP	TP	FN	AUC-ROC	Recall
GradientBoost	1916	418	89	34	0.861	0.72
AdaBoost	2020	314	80	43	0.869	0.65
XG Boost	1827	507	94	29	0.864	0.76
Voting Classifier	1902	432	91	32	0.863	0.74

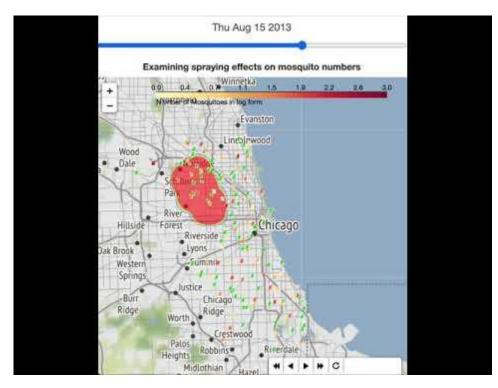
- L. INTRO & EDA Mak
- PREPROCESSING,
- 2. **FEATURE SELECTION & MODELING**Alex
- 3. MODEL EVALUATION Robby
- 4. **COST BENEFIT ANALYSIS** Elliot



# Cost-Benefit Analysis

# Overview of Spraying Control Measures

- Currently no effective treatments for WNV
- Prevention of disease typically relies on pest management programmes and control measures (e.g. spraying)
- Due to lack of relevant trap data immediately before and after spraying, cannot conclude that spraying is able to reduce the number of positive WNV cases
- However, several studies have illustrated effectiveness of spraying as both a reactive/proactive control measure to combat WNV



# Reactionary Control Measures

- Taken after a WNV outbreak
- Tend to be associated with **exorbitantly high costs** (typically in millions USD) due to the high intensity needed

### 1. Sacramento-Yolo Mosquito and Vector Control District (2005)

- a. Application of adulticide indicated a 75% reduction of *Culex pipiens* and a 48.7% reduction in *Culex tarsalis* population in the treated area
- b. Human WNV case incidence was significantly lower in the treated area

### 2. Texas WNV outbreak (2012)

a. WNV neuroinvasive disease incidence decreased from 7.31/100,000 before treatment to 0.28/100,000 after treatment in the treated area

### 3. Chicago WNV outbreak (2005)

- a. Reduction of adult mosquito abundance by 54% in the treated areas.
- b. During this same period, mosquito abundance actually increased by 153% in the untreated areas

### Proactive Control Measures

• Taken prior to any WNV outbreak, usually as part of a surveillance system:

#### 1. Atlanta

a. Larvicides was applied to catch basins in urban park areas over the course of 2 seasons, resulting in more than 90% decrease in larval/pupal production in the catch basins

#### 2. **Southern California**

- Limited reactive aerial spraying was noted to be insufficient to reduce vector mosquito abundance, WNV infection rates in mosquitoes or spread of WNV transmission
- b. When intensive spraying was done at the first detection of WNV in mosquitoes, average of 61% reduction in vector abundance and also associated reductions in WNV infection rates

### 3. **Sacramento (2007)**

- a. WNV infection rates in *Culex tarsalis* and *pipiens* mosquitoes had exceeded levels of concern established
- b. Spraying of a 215 sq km area each day for 3 consecutive days resulted in a 57% and 40% reduction in *tarsalis* and *pipiens* abundance respectively, and the WNV infection rate declined by 77% and 21% in *tarsalis* and *pipiens* mosquitoes respectively during the following 3 days

# Reactionary vs. Proactive approach

- Both reactive and proactive spraying are useful in controlling the WNV carrier mosquito population, and WNV infection rates
- Reactive efforts are far more expensive in terms of cost of control measures, as well as economic costs
  resulting from human infections already occurred
- Proactive efforts include establishing of surveillance systems for monitoring mosquito numbers and WNV infection rates, and appropriate response plans when particular threshold levels are met
  - Factors could include counts of dead birds, temperature, rainfall
  - Coupled with localized knowledge of vector hotspots and experience with previous outbreaks
  - Typically more effective and less costly
- Our model is an example of the proactive approach
  - Predict locations that are likely to have WNV occurrences based on weather data and historical factors
  - Perform targeted spraying of specific areas during months when virus-transmission rates are high
  - Determine extent and frequency of spraying to be performed

# Historical Costs of Vector Control

- 1. **New York State, city and 4 counties (1999)**: The state, city and 4 counties spent more than USD 14 million on protective measures such as mosquito control.
- 2. **St. Tammany Parish (Louisiana 2002)**: Additional mosquito control activities during the 2002 WNV outbreak cost USD 1.7 million over their usual USD 2 million budget.
- 3. **Sacremento (2005)**: The district spent USD 700,000 on aerial ULV applications alone in response to the 2005 WNV outbreak.
- 4. **Dallas County (2012)**: Similarly, aerial ULV applications during the 2012 outbreak cost approximately USD 1.7 million.

It is worth noting that in these 4 instances the measures were reactionary, typically in response to an occurring outbreak and that the costs were exorbitantly high just to control it.

### **Annual Costs**

- 1. Contractual price for spraying
- Damage to non-target organisms, humans, environments due to spraying

#### **Assumptions:**

- Chicago's land area is 601.6 km²
- Unit cost of spraying on truck mounted ULV: USD 35/km
- Unit cost of spraying using backpack ULV: USD 560/km
- Chicago has more than 6,450 km of streets and 3,050 km of alleys. To fully spray the whole area of Chicago,
  - Streets will be fully used
  - Alleys required are 40% of the total alley length
- Spraying will be applied at least 4 times during the mosquito season (Aug/Sep) as a one-time application lasts about 14 -21 days
- Ecological and human impact of spraying are taken to be 0



### Total costs:

~ 4 million USD/year

# Annual Benefits

- 1. Avoidance of medical treatment costs
- 2. Avoidance of productivity loss

### **Assumptions:**

- Average number of cases of WNV per year (from 2012 to 2018) is 35, considering spraying regime in place
  - Should there be no spraying, there could be  $\sim 35 \times 6 = 210 \text{ cases/year}$
  - Spraying helps to prevent ~ 175 cases/year
- Overall median cost due to treatment and loss of productivity is USD 28,000

F	Referenced to 2012 USD	Fever (N=18)	Meningitis (N=19)	Encephalitis (N=16)	AFP (N=27)
	Direct costs:				
	Inpatient hospital costs	4,467 (419-23,374)	7,261 (337-13,633)	15,136 (3,734-207,303)	20,774 (5,066-264,176)
	Lost productivity	328 (92-2,729)	682 (68-1,592)	1,380 (113-307,871)	2,136 (232-145,750)
	Total direct initial costs	4,617 (538-24,010)	7,942 (1,057-14,569)	20,105 (3,965-324,167)	25,117 (5,385 - 283,381)
	Long-term costs:	Fever (N=12)	Meningitis (N=11)	Encephalitis (N=5)	AFP (N=10)
	Medical appointments	109 (0-677)	0 (0-851)	495 (0-17,160)	3,671 (452-12,093)
	Additional care costs	0 (0-8,900)	0 (0)	334 (0-10,013)	278 (0-6,119)
	Medicines & equipment	72 (0-5,320)	33 (0-1,305)	109 (0-1,964)	590 (106-427,028)
	Lost productivity	1,180 (0-39,760)	10,363 (0-258,592)	0 (0-5,596)	6,771 (0-143,033)
	Total long-term costs	2,271 (0-41,401)	10,556 (0-260,748)	8,055 (0-23,693)	22,628 (624-439,945)
Simple overal	ll median (not from paper)	6,888	18,498	28,160	47,745

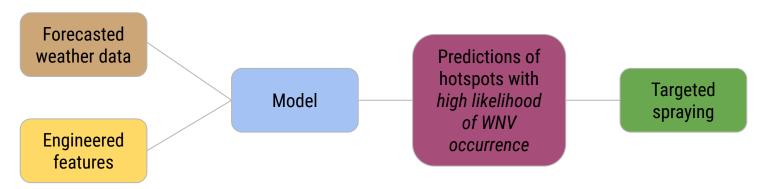
Associated costs (median with range) due to WNV complications

### Total benefits:

~ 5 million USD/year

### Recommendations

- The net benefit of the current vector management program in Chicago is hence approx. **USD 1 million/year**
- However, the assumed USD 4 million/year spent on spraying is unsustainable, given Chicago's limited public health budget of USD 36 million/year



- Reduce cost of spraying by ~30% initially (USD 1.2 million savings/year)
- Assume that our targeted spraying also drives average infection rates down by 50% (USD 500,000 savings/year)
  - Ultimately, our model is able to deliver annual benefits of USD 1.7 million/year
- We also recommend the adoption of a complete, proactive vector management plan

# Conclusion

### Conclusion

- Our VotingClassifier model is fully predictive (does not require the wnv\_species feature, which cannot be forecasted before hand)
  - Ensemble of GradientBoostingClassifier, AdaBoostClassifier and XGBClassifier boosting models
  - Achieved a validation AUC of **0.852** and recall of **0.74** on the validation dataset
  - Kaggle submission achieved an AUC score of 0.718
- Model is able to aid the city of Chicago's administration in determining the hotspot locations
  - Perform targeted spraying of pesticides as a proactive form of mosquito control and reduce costs to manageable level
  - Protect its people from being afflicted by the potentially life threatening disease
- Future improvements to our model:
  - Incorporate other potentially relevant features such as counts of dead birds, proximity to water bodies
  - Obtaining more granular data to better account for the efficacy of spraying, e.g.
    - Regular daily trap testing data for a week prior and after spraying
    - Data for all traps within the spray areas
    - Data for sprayed areas that have historically high number of mosquitos to better observe the spray effect

# Thank you!



# Models comparison

	TN	FP	TP	FN	AUC-ROC	Recall
GradientBoost	1916	418	89	34	0.861	0.72
AdaBoost	2020	314	80	43	0.869	0.65
XG Boost	1827	507	94	29	0.864	0.76
Voting Classifier	1902	432	91	32	0.863	0.74
Voting exclu. species	1884	450	91	32	0.857	0.74

### Adult Mosquito Fogging Pricing (City of Burleson - Texas, 2016)

#### Source

#### ATTACHMENT "A"

#### Adult Mosquito Fogging Pricing Schedule

#### Adult Mosquito Control Services as Needed

Truck mounted ULV Cost per linear mile: \$40.00 minimum application rate

\$46.00 mid-level application rate

UTV mounted ULV Cost per 1/2 linear mile:

Paved Surfaces (cart paths) \$160.00 mid-level application rate

Off Road Area \$200.00 mid-level application rate

Back-pack ULV Cost per 1/4 linear mile: \$190.00 mid or maximum level application rate

Submission and Description	Private Score	Public Score	Submission and Description	Private Score	Public Score
voting_model.csv just now by Robby Sim voting model final	0.75110	0.75690	voting_model_drop_species.csv just now by Robby Sim Voting_model_drop_species	0.71753	0.71890
xgb_model.csv a minute ago by Robby Sim xgb model final	0.69293	0.69345			
gb_model.csv 2 minutes ago by Robby Sim gb_model final	0.71155	0.72006			
ab_model.csv 2 minutes ago by Robby Sim ab_model final	0.70631	0.71298			