

# Verifying Volunteer Entries to the NYC Trees Census

Joshua Szymanowski



# A look at the data



**NYC** OpenData NYC Parks

TREES  
COUNT 2015

**652,173** trees

444,390 professional (TreesCount and Parks staff) entries

207,776 volunteer entries

**Health** of trees (classification target):

**Good** (82%), **Fair** (14%), and **Poor** (4%)

**Variables:**

Tree diameter, species, number of stewards, quality of tree guards,  
root/trunk/branch problems

**Location attributes:**

Borough, neighborhood, community board, council district, state  
assembly, state senate

Latitude & longitude for **mapmaking**

# Project Outline

1

2

3

4

5

## Data Cleaning

Drop stumps and dead trees  
Drop volunteer data  
Turn data objects into  
numericals

## EDA

Effects of location:  
Borough  
Neighborhood  
Political districts  
Effect of species  
Root, trunk, and branch  
problems  
Tree guards  
Tree stewards  
Maps

## Baseline Model

Fit a Random Forest Model  
with standard parameters  
  
Interpret which features are  
most important and if any  
more EDA needs to be done

## Feature Engineering

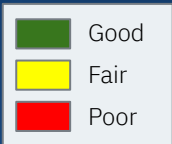
Distance of nearest tree  
using GeoPandas  
  
Number of trees on same  
block  
  
Create dummy variables  
  
Use community board as  
neighborhood variable

## Final Model

Random Forest with  
GridSearch  
  
Naive Bayes with  
GridSearch  
  
Create list of flagged  
trees (future goal)



## Volunteer entries



## Professional entries





# *Goals*

Recommend changes for next tree census

Recommend policies for trees planting and maintenance

Develop a model to verify the health status given by volunteers



# Model preview

## Vanilla random forest

Hyperparams: *class\_weight='balanced'*

Accuracy: **54.7%**

Weighted F1: **61.5%**

## Untuned random forest

Hyperparams: *class\_weight='balanced'*

Accuracy: **82.1%**

Weighted F1: **77.8%**

## Tuned random forest

Accuracy: **72.1%**

Weighted F1: **74.5%**

Hyperparams:

*max\_features=11,*

*max\_depth=55,*

*min\_samples\_leaf=2,*

*n\_estimators=500,*

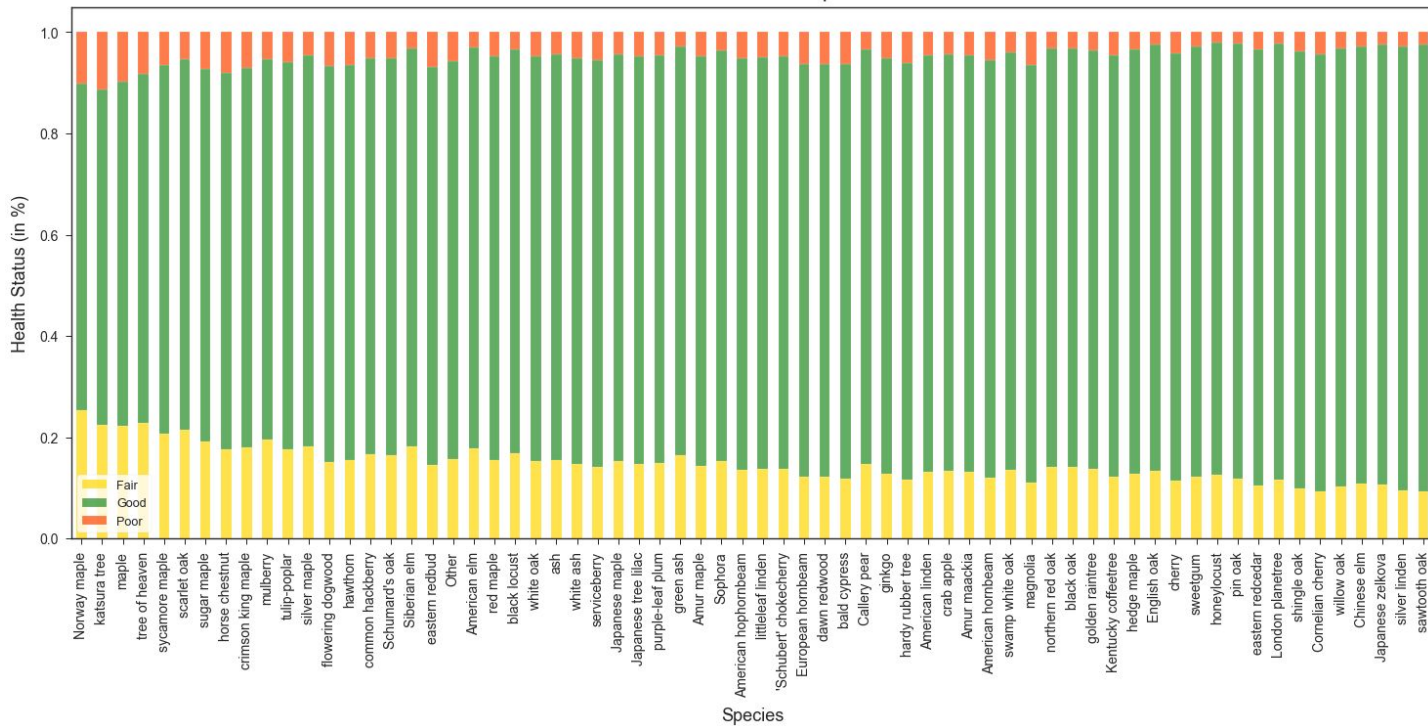
*class\_weight='balanced'*

# Species matters

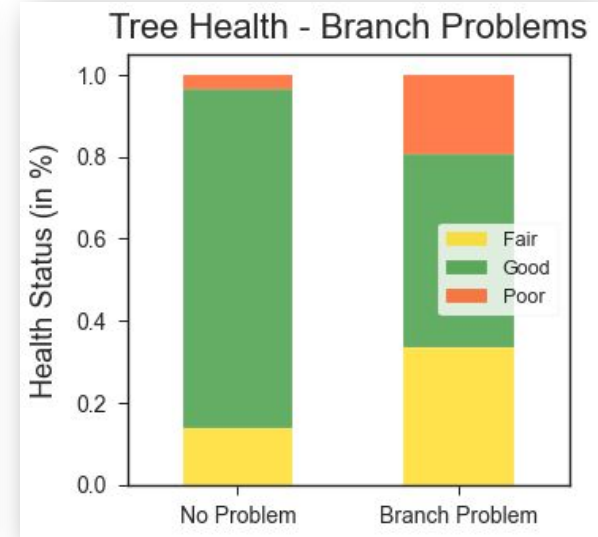
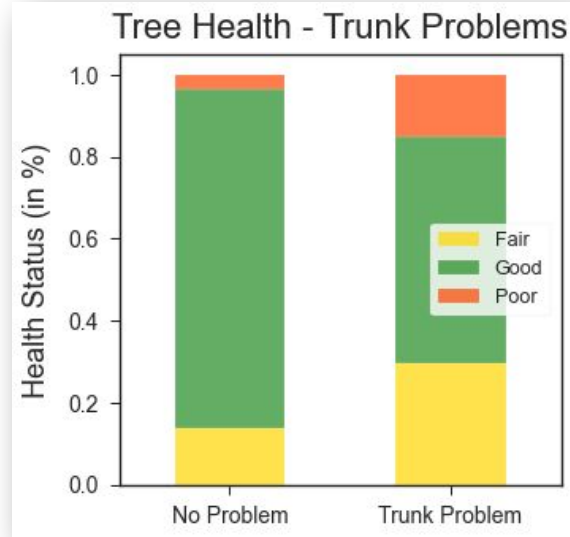
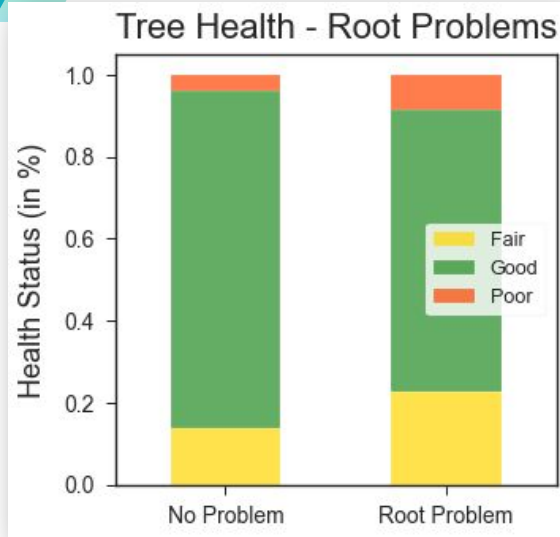
**Top 3:**  
Norway maple  
Katsura tree  
Maple

**Bottom 3:**  
Sawtooth oak  
Silver linden  
Japanese zelkova

Tree Health Across Species



# Tree problems (and solutions)



Solutions for healthier trees:

The most important problems are listed as "Other".

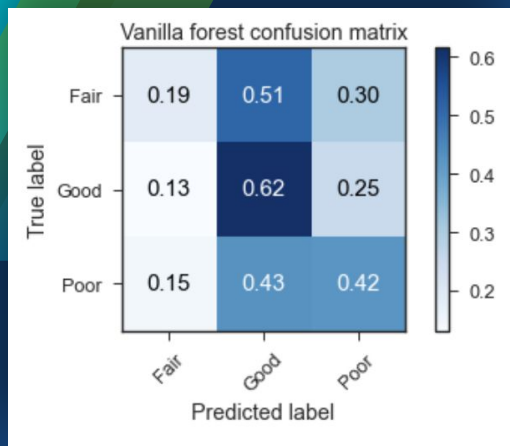
In future censuses, be more specific or have a notes column for each.

More regular maintenance of trees.

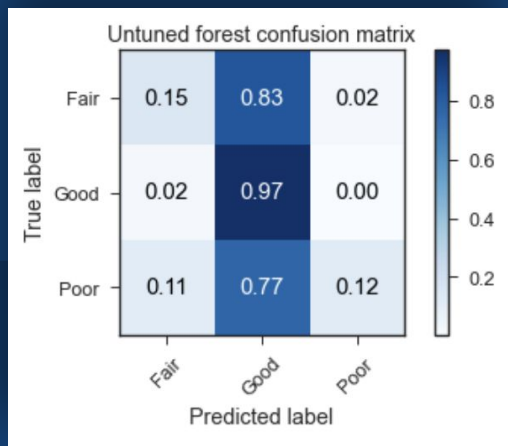
Farther reaching environmental protections, à la the plastic bag ban.



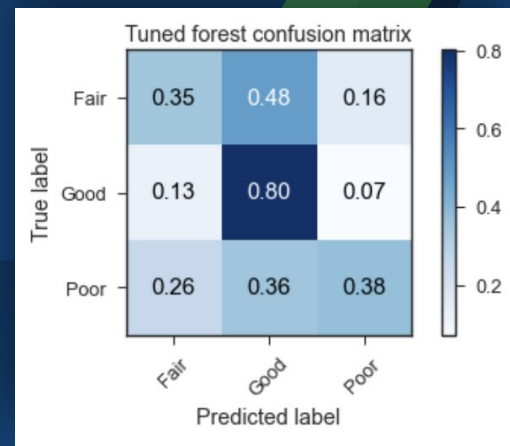
# Confusion journey



Accuracy: **54.7%**  
Weighted F1: **61.5%**



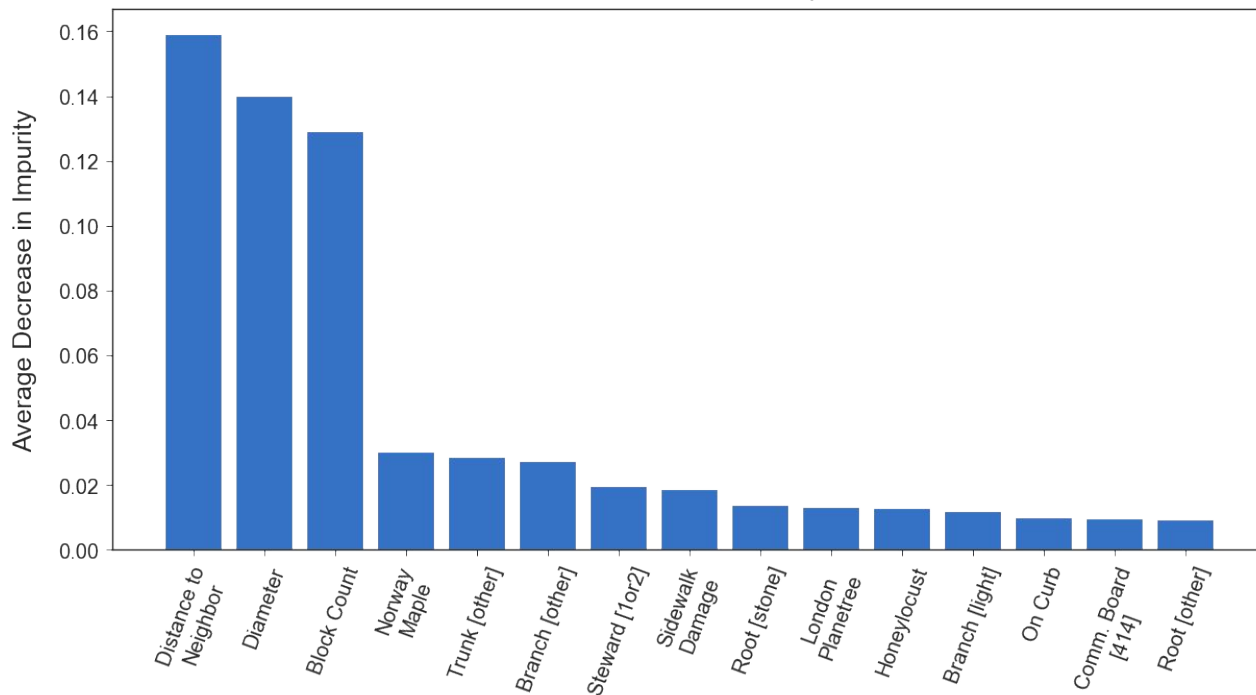
Accuracy: **82.1%**  
Weighted F1: **77.8%**



Accuracy: **72.1%**  
Weighted F1: **74.5%**

# Final model - Random Forest

Random Forest Feature Importances



Accuracy: **72.1%**

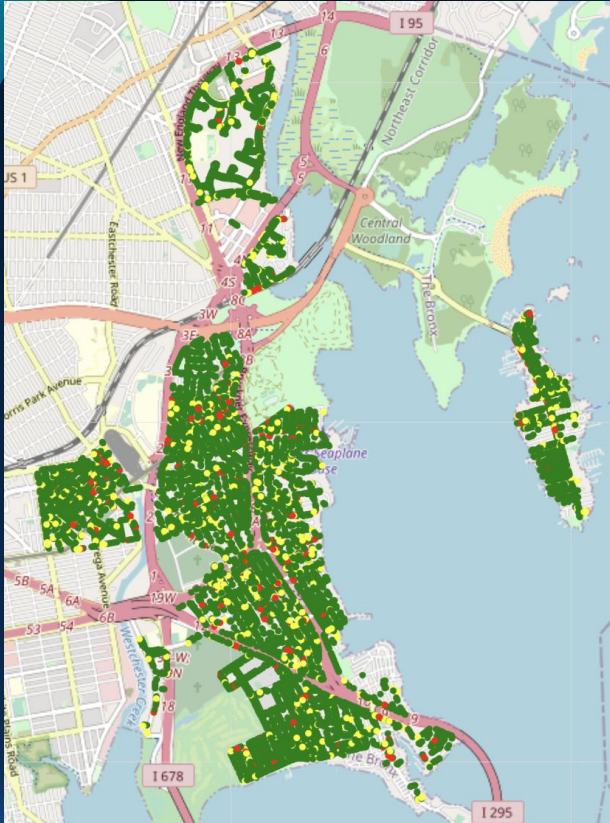
Weighted F1: **74.5%**

Weighted Precision: **72.1%**

Weighted Recall: **77.6%**

## Community Board 210

Neighborhoods: *Co-op City, City Island, Throggs Neck, Country Club, Zerega, Westchester Square, Pelham Bay, Waterbury Lasalle*



## Community Board 414

Neighborhoods: *Breezy Point, Belle Harbor, Broad Channel, Neponsit, Arverne, Bayswater, Edgemere, Rockaway Park, Rockaway and Far Rockaway*

# Any questions?