

# Project Jetson

---

Katherine Hoffmann Pham

August 15, 2019

Data Fellow, UN Global Pulse

Ph.D. Candidate in Information Systems, NYU Stern School of Business

# Overview

---

# About UN Global Pulse

- **UN Global Pulse:** “innovation initiative of the UN Secretary-General . . . to accelerate discovery, development and scaled adoption of big data innovation for sustainable development and humanitarian action”
- **Model:** partner with other UN entities on projects to demonstrate the value of big data



# About UN Global Pulse

- Motivation:
  - SDGs: Timely, cost-effective monitoring of 169 targets
  - Appeal of “always-on”, “real-time” big data sources
- Sample projects:
  - Measuring Poverty With Machine Roof Counting
  - Making Ugandan Community Radio Machine-Readable



# About the Fellows Program

## Global Pulse Data Fellows Program:

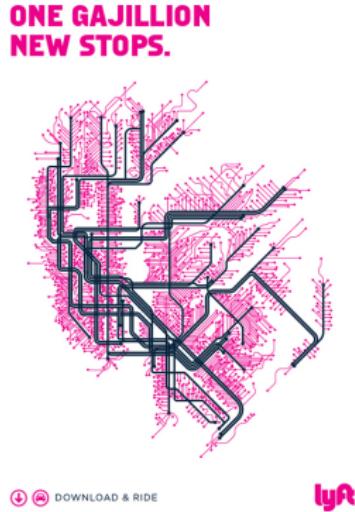
Match PhD students with UN entities for 6-12 months

- 40+ projects applied from within UN system
- 8 projects selected from first cohort
- Example projects:
  - Parsing radio broadcasts for disease surveillance (WHO)
  - Identifying refugee settlements in satellite data (UNOSAT)
  - Exploring the evolution of #metoo movement on social media

# About Me

## Thesis: Digital Future of Mobility

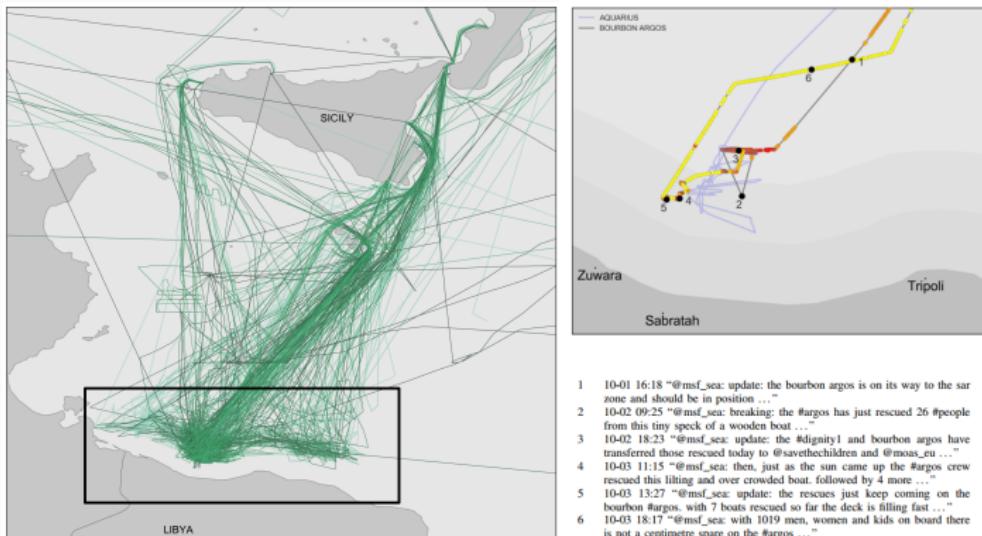
- Impact of Uber/Lyft on public transportation in NYC
- Nudges to improve parking among e-scooter users



# About Me

## Projects: Big Data for Social Good

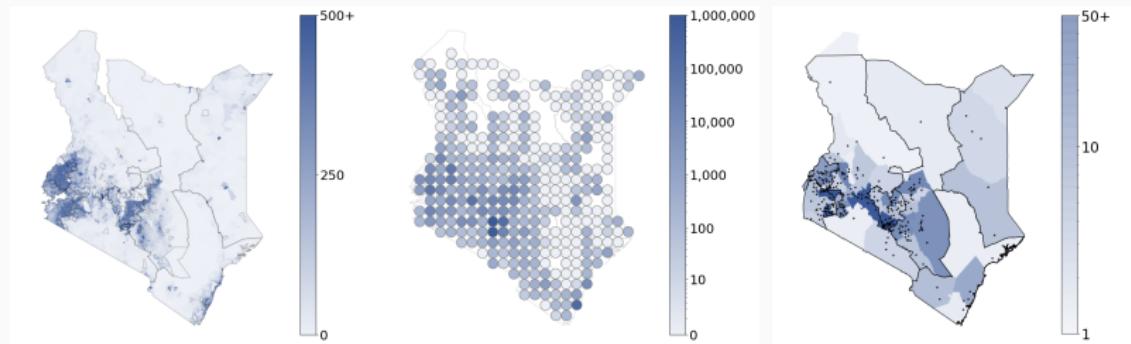
- Migration in the Central Mediterranean using AIS data



# About Me

## Projects: Big Data for Social Good

- Facebook advertising data for demographic estimates in Kenya



# About Me

Training: BIGSSS Summer School on Migration

- How do people model movement?
  - Agent-Based Models (rules + simulation)
  - Gravity Models (flows between two regions)
- What is the focus of current research?
  - Host community based (integration, political orientation)
  - Theoretical (climate change, resilience, and agency)
  - Global models (country-level flows)

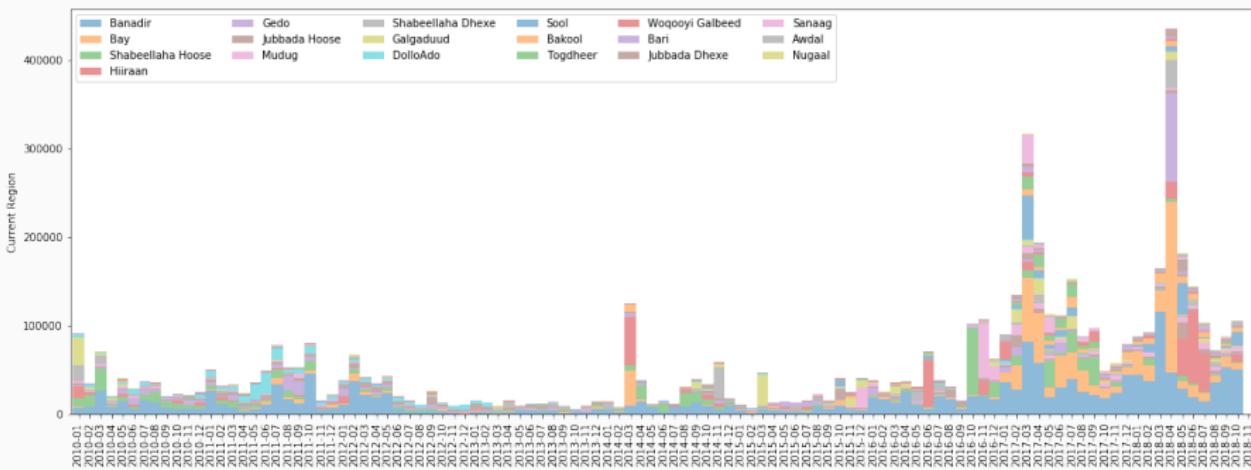
# Project Jetson

---

# Project Jetson

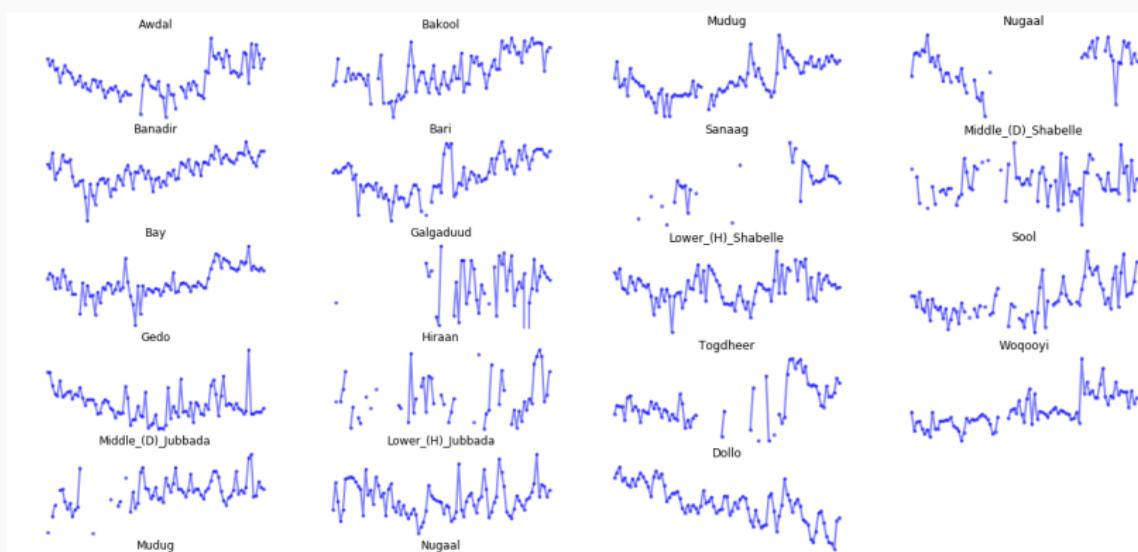
## Research Questions

Can we use artificial intelligence to forecast the monthly volume of IDP/refugee arrivals by region in Somalia?



# Research Questions

Can we use artificial intelligence to forecast the monthly volume of IDP/refugee arrivals by region in Somalia?



# Research Questions

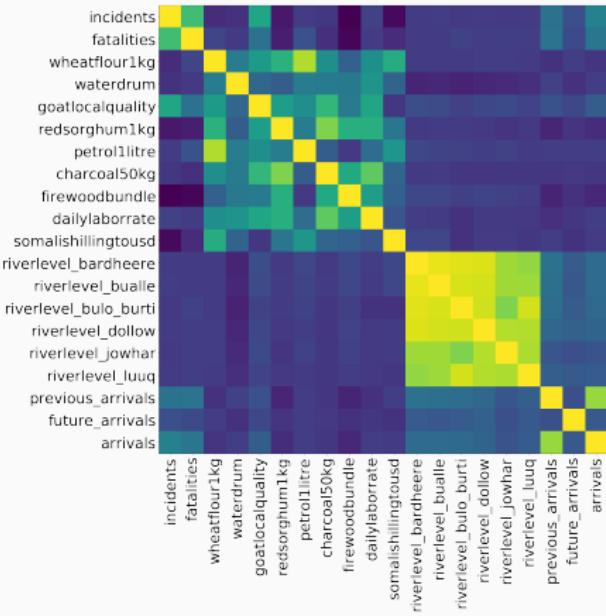
Can we use artificial intelligence to forecast the **monthly volume of IDP/refugee arrivals by region** in Somalia?

- How far in advance:  
1 months? 3 months?
- How does **contextual data** improve predictions?
- How best to **model** this problem: arrivals? flows?  
changes?

# Research Questions

Can we use artificial intelligence to forecast the monthly volume of IDP/refugee arrivals by region in Somalia?

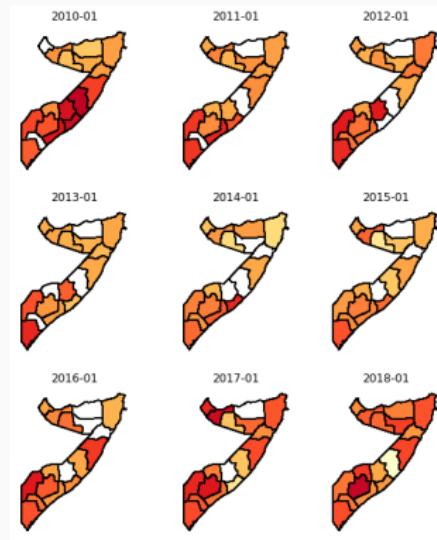
- How far in advance:  
1 months? 3 months?
- How does contextual data improve predictions?
- How best to model this problem:
  - arrivals?
  - flows?
  - changes?



# Research Questions

Can we use artificial intelligence to forecast the monthly volume of IDP/refugee arrivals by region in Somalia?

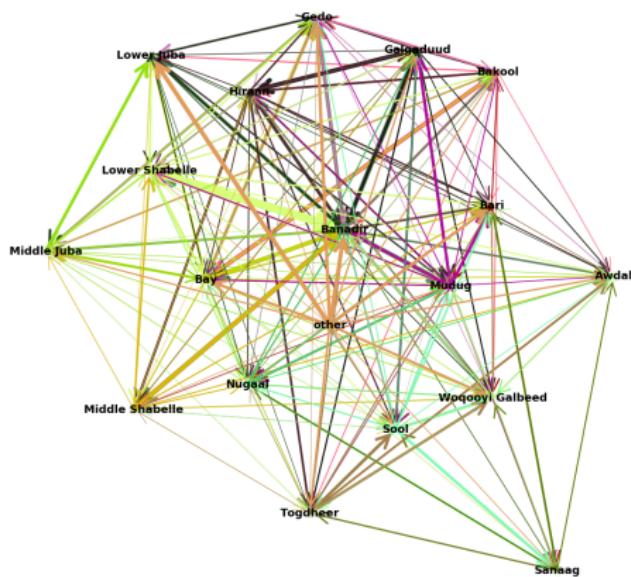
- How far in advance:  
1 months? 3 months?
- How does contextual data improve predictions?
- How best to model this problem: arrivals? flows? changes?



# Research Questions

Can we use artificial intelligence to forecast the monthly volume of IDP/refugee arrivals by region in Somalia?

- How far in advance:  
1 months? 3 months?
- How does contextual data improve predictions?
- How best to model this problem: arrivals? flows? changes?

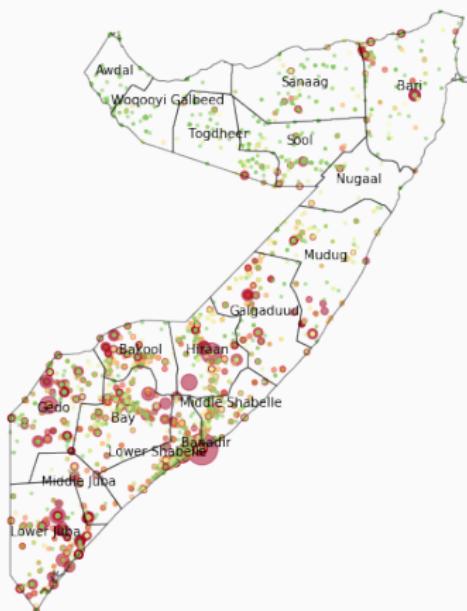


# Experiment 1

- Data extracts are gathered and compiled monthly  
*Can this be automated?*
- $18 + 1$  regions  $\times$  23 models per region  
*Can we implement a single model across all regions?*
- Goal is to predict arrivals one month in advance  
*Can we predict 3 months in advance?*
- Models implemented in Eureqa (equations)  
*Can we explore alternative tools?*
- Each month, 3 winners are chosen based on previous month  
*Can this be automated?*
- Models are evaluated using several characteristics  
*Can we identify the most important metric for the team?*

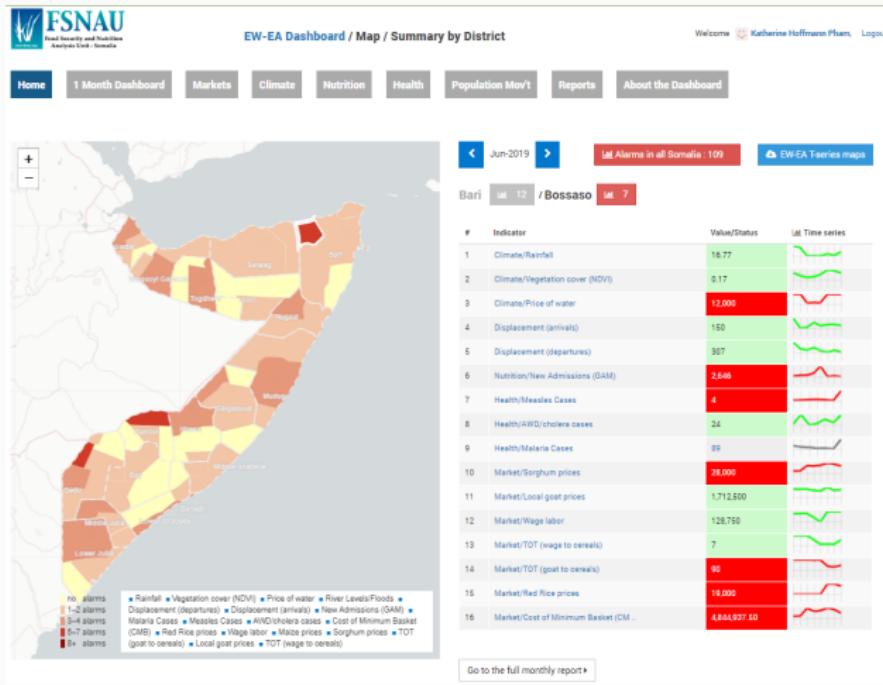
# Data Sources

ACLED: conflict incidents and fatalities



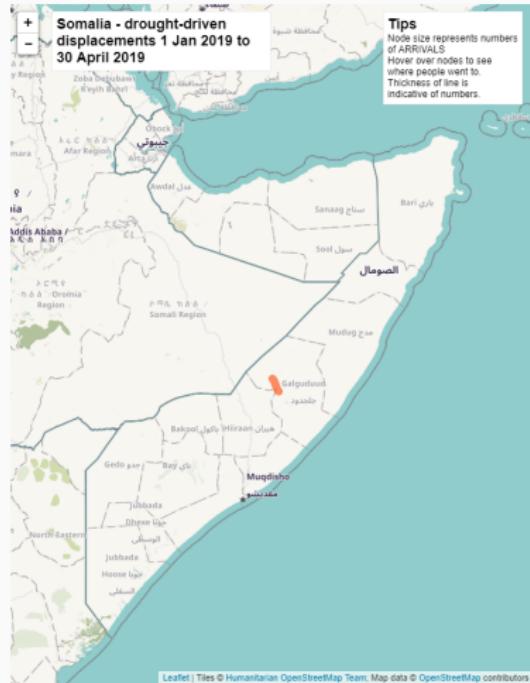
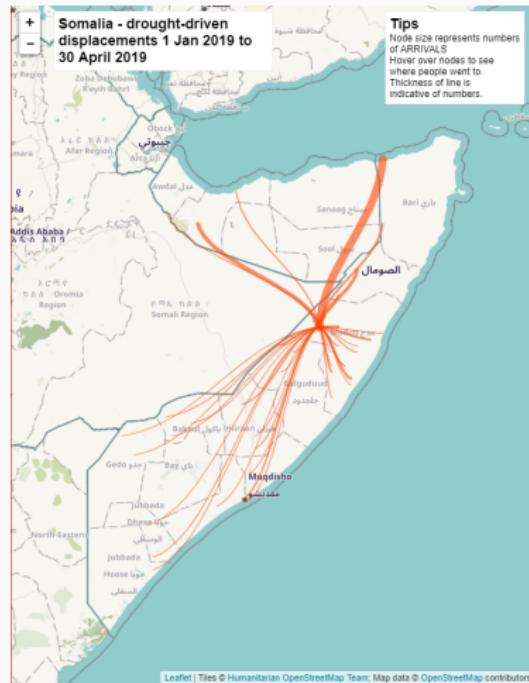
# Data Sources

## FSNAU Early Warning Early Action Dashboard: prices; climate; health



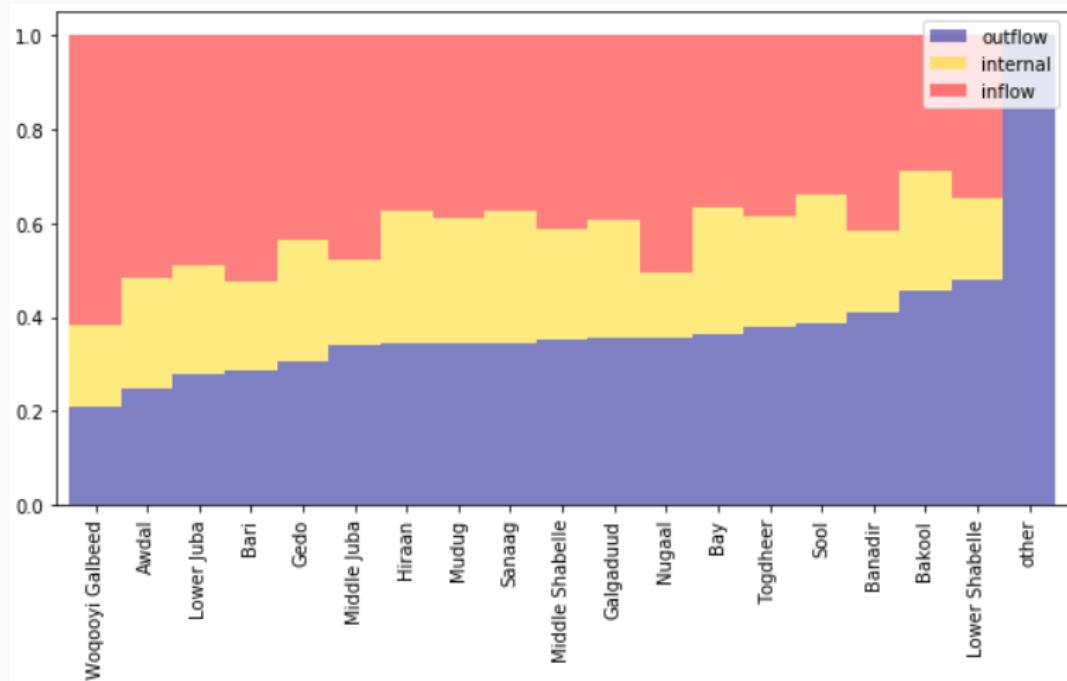
# Data Sources

PRMN: current, previous, or future location



# Data Sources

PRMN: current, previous, or future location



# Problem Setup

1. Model output: # monthly arrivals in region
2. Model input:
  - Lagged data for the region  
(e.g. previous 3,4,5,6,12 months)
  - Lagged data for other regions  
(e.g. previous 3,4,5,6,12 months)
  - Indicator for missing values
  - Indicator for region and month of year
  - Linear time trend

# Tools tested

## Python (SciKitLearn)

```
# ...

# Specify model
m = DecisionTreeRegressor(min_samples_leaf=5, max_depth=5, criterion='mse')

# Fit model
m.fit(X_train, y_train)

# Make predictions
y_fit = m.predict(X_train)
y_pred = m.predict(X_test)

# ...
```

# Tools tested

H2o.ai

H2o.ai Experiment tuhenuma

DATALESS AI 15.4 - AI TO DO AI  
Licensed to University of New York (SNC04073 - Academic License). Current User - XH0F212

ASSISTANT

STATUS: COMPLETE

EXPERIMENT SETTINGS

EXPERT SETTINGS

SCORER

- GINI
- R2
- MSE
- RMSE**
- RMSLE
- RMSPE
- MAE
- MER
- MAPE
- SMAPE

CPU / MEMORY

Notifications Log Trace

ITERATION DATA - VALIDATION

VARIABLE IMPORTANCE

ACTUAL VS PREDICTED

SUMMARY

Experiment: tuhenuma, 2019-06-07 15:54, 15.4  
Setting: 10/7/70, 24/48 -476337796, GPUs disabled  
Train data: learn\_df\_region\_20190607\_train (1615, 2565)  
Validation data: N/A  
Test data: learn\_df\_region\_20190607\_test (108, 2565)  
Target column: arrivals (regression)  
System specs: Linux, 12 GB, 4 CPU cores, 0/0 GPU  
Max memory usage: 5.49 GB, 0 GB GPU  
Recipe: AutoDL (227 iterations, 8 individuals)  
Validation scheme: time-based, 4 internal holdouts  
Feature engineering: 8260 features scored (10 selected)  
Timing:  
Data preparation: 116.57 secs  
Model and feature tuning: 4132.17 secs (305 of 780 models trained)  
Feature evolution: 6900.60 secs (2400 of 2432 models trained)  
Python / MOJO scoring pipeline building: 32.98 secs / 0.00 secs  
Validation score: RMSE = 0329.5 +/- 2913.9 (baseline)  
Validation score: RMSE = 64081.1 +/- 39598.7 (final pipeline)  
Test score: RMSE = 3941.4 +/- 523.03 (final pipeline)

The screenshot shows the H2o.ai Experiment tuhenuma interface. At the top, it displays basic experiment statistics: 2K rows, 3K columns, 0 dropped cols, and a validation dataset set to 'Yes'. Below this, the 'Dataset' section shows 'learn\_df\_region\_20190607\_train' with target column 'arrivals' and type 'real'. The 'Assistant' section indicates the status is 'COMPLETE'. The 'Experiment Settings' section has values for Accuracy (10), Time (7), and Interpretability (10). The 'Expert Settings' section lists various scorers: GINI, R2, MSE, RMSE (selected), RMSLE, RMSPE, MAE, MER, MAPE, and Smape. The 'CPU / Memory' section shows current usage. The 'Iteration Data - Validation' section includes a histogram of iteration data from 0 to 240, with a peak around 200. The 'Variable Importance' section lists variables like '2541\_TargetLogregion.9' with a value of 1.00. The 'Actual vs Predicted' section shows a scatter plot comparing actual vs predicted values. The 'Summary' section provides detailed timing and system specifications.

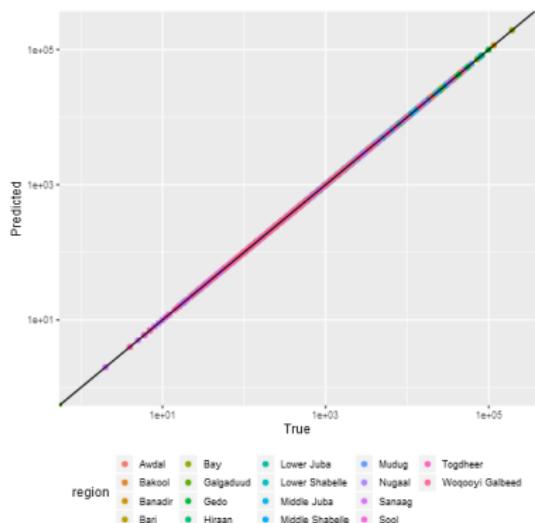
# Models tested

## Ground truth:

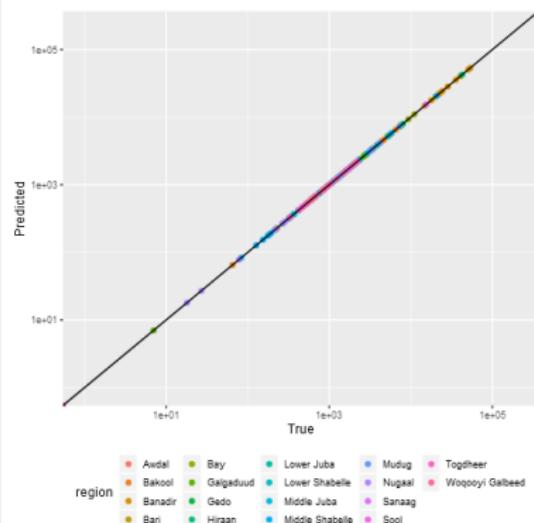
Choose model:

true

Train dataset



Test dataset



# Models tested

Baseline: Last observation carried forward (1 month)



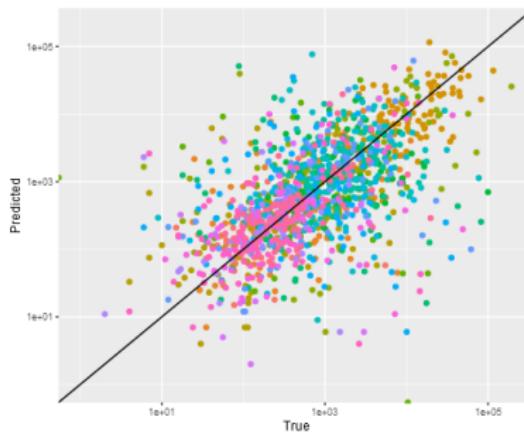
# Models tested

Baseline: Last observation carried forward (3 months)

Choose model:

locf3

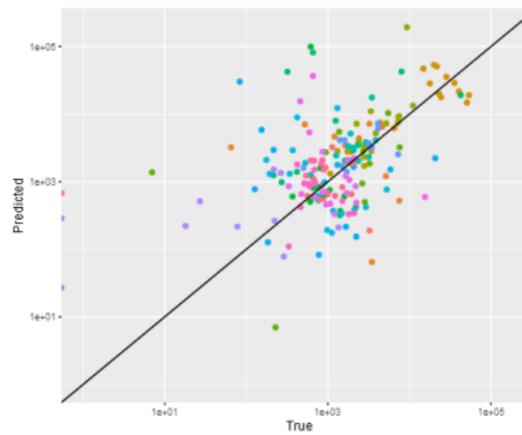
Train dataset



region

• Awdal  
• Bakool  
• Banadir  
• Bari  
• Gedo  
• Hirran  
• Jubbada Hoose  
• Jubbada Dhexe  
• Lower Shabelle  
• Middle Shabelle  
• Nugaal  
• Sanag  
• Sool  
• Togheer  
• Woqooyi Galbeed

Test dataset



region

• Awdal  
• Bakool  
• Banadir  
• Bari  
• Gedo  
• Hirran  
• Jubbada Hoose  
• Jubbada Dhexe  
• Lower Shabelle  
• Middle Shabelle  
• Nugaal  
• Sanag  
• Sool  
• Togheer  
• Woqooyi Galbeed

# Models tested

Linear regression:  $\approx$  try to fit a line to the data



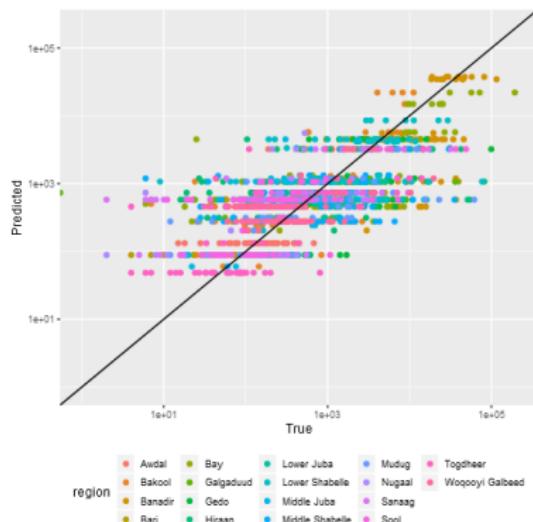
# Models tested

Decision trees:  $\approx$  break the dataset into groups

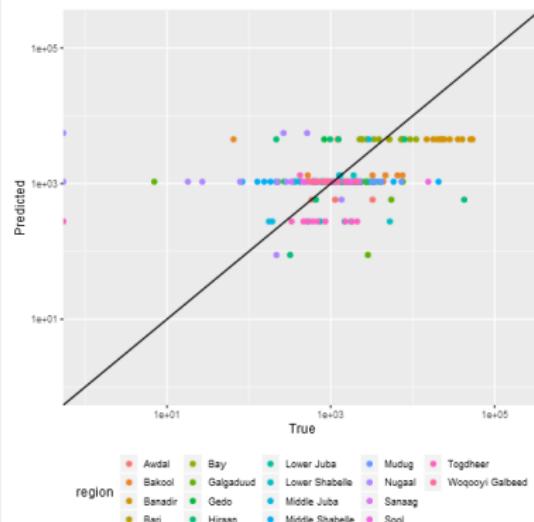
Choose model:

Decisiontree

Train dataset



Test dataset



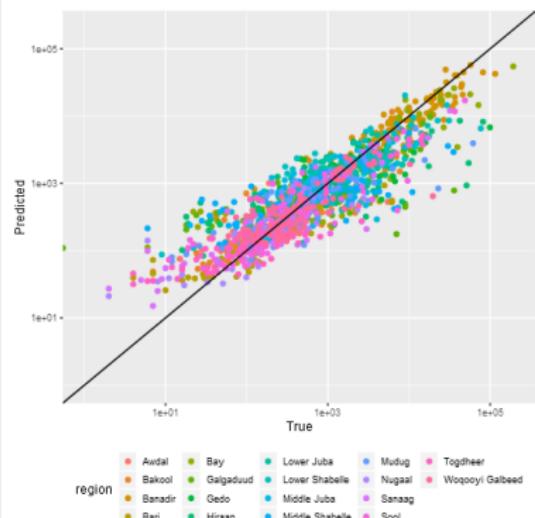
# Models tested

Boosting: combine many models in a smart way

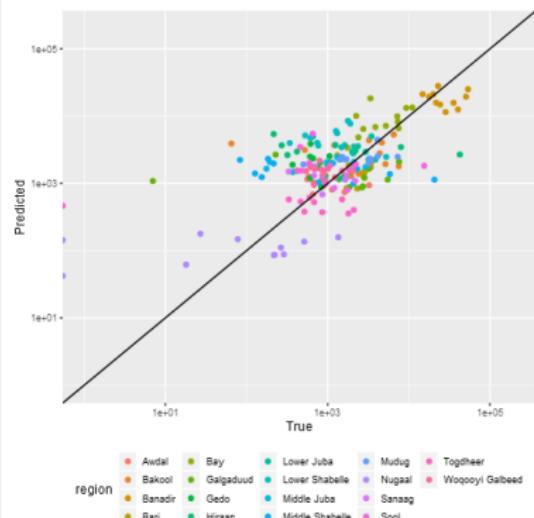
Choose model:

Xgboost

Train dataset



Test dataset



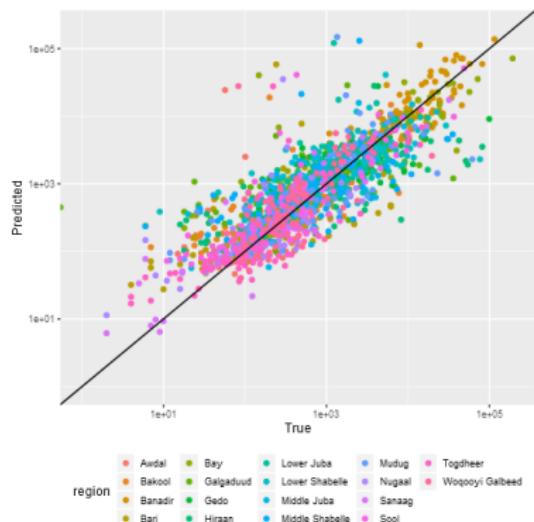
# Models tested

## Perceptron: fit a very flexible model

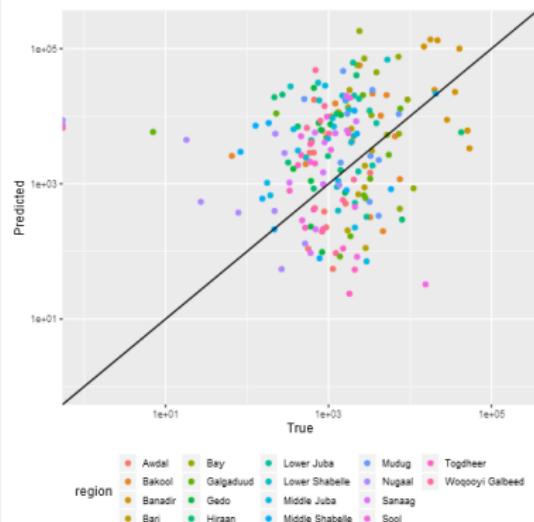
Choose model:

Perceptron

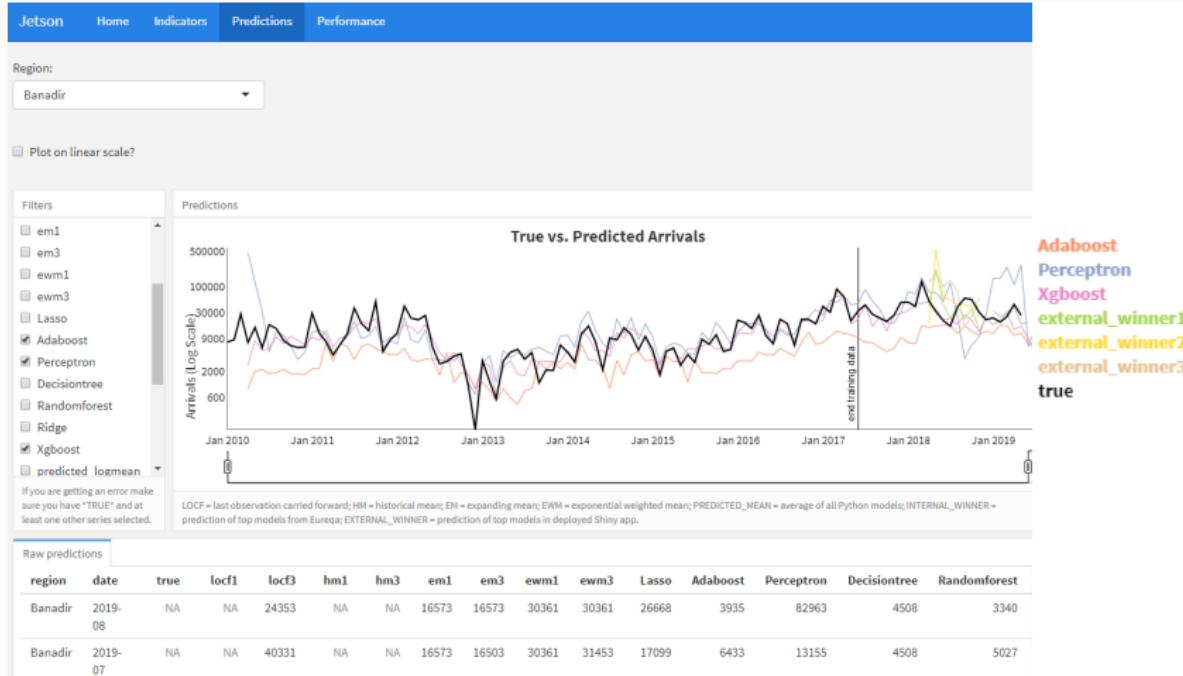
Train dataset



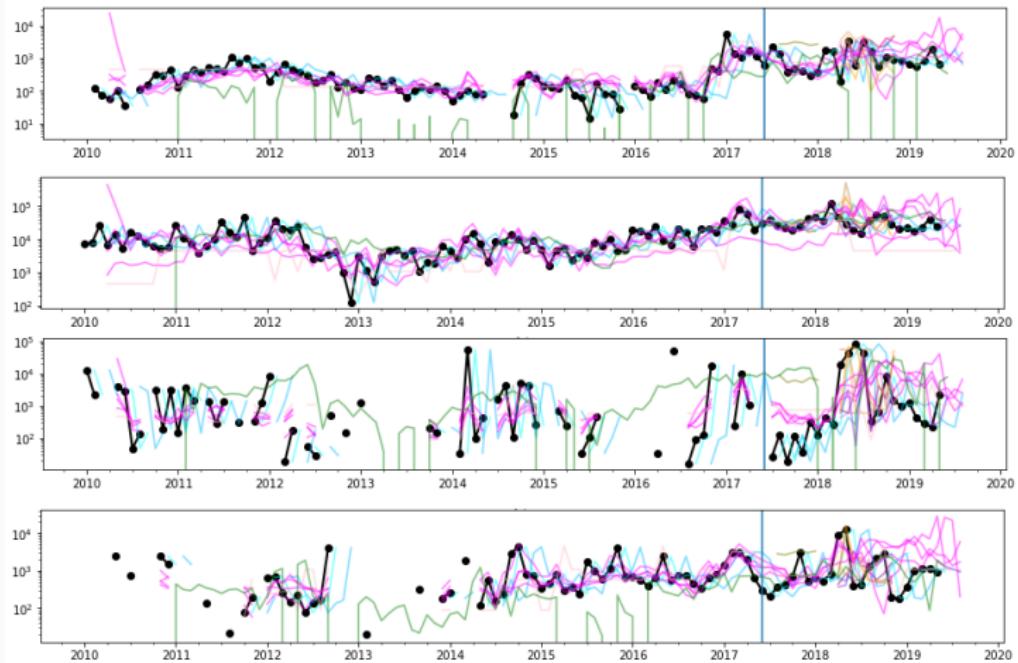
Test dataset



# Model Comparison



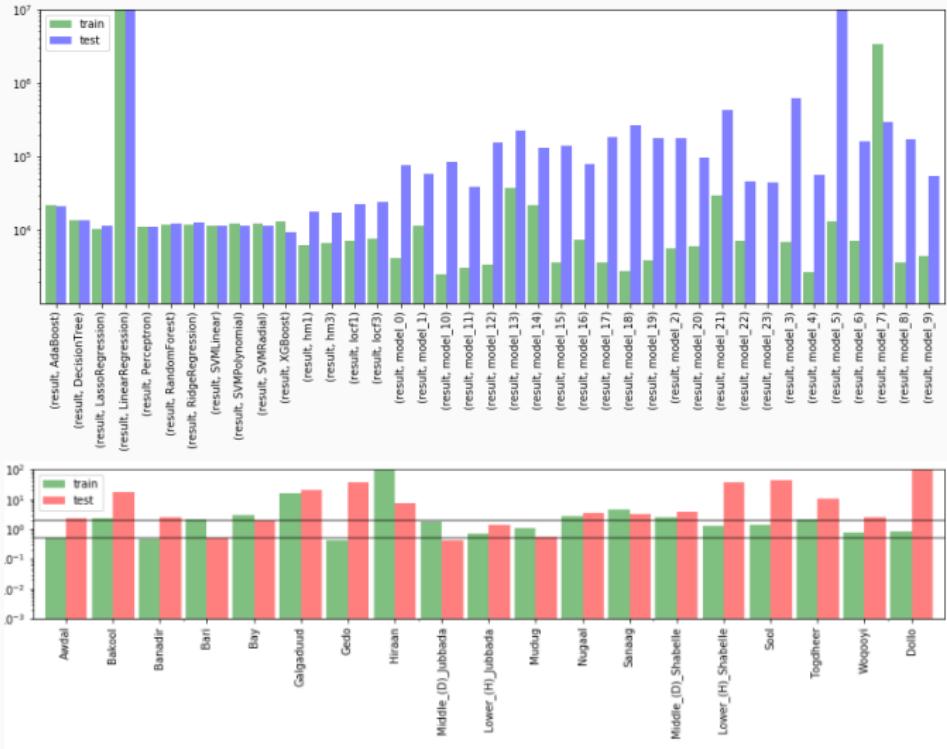
# Model Comparison



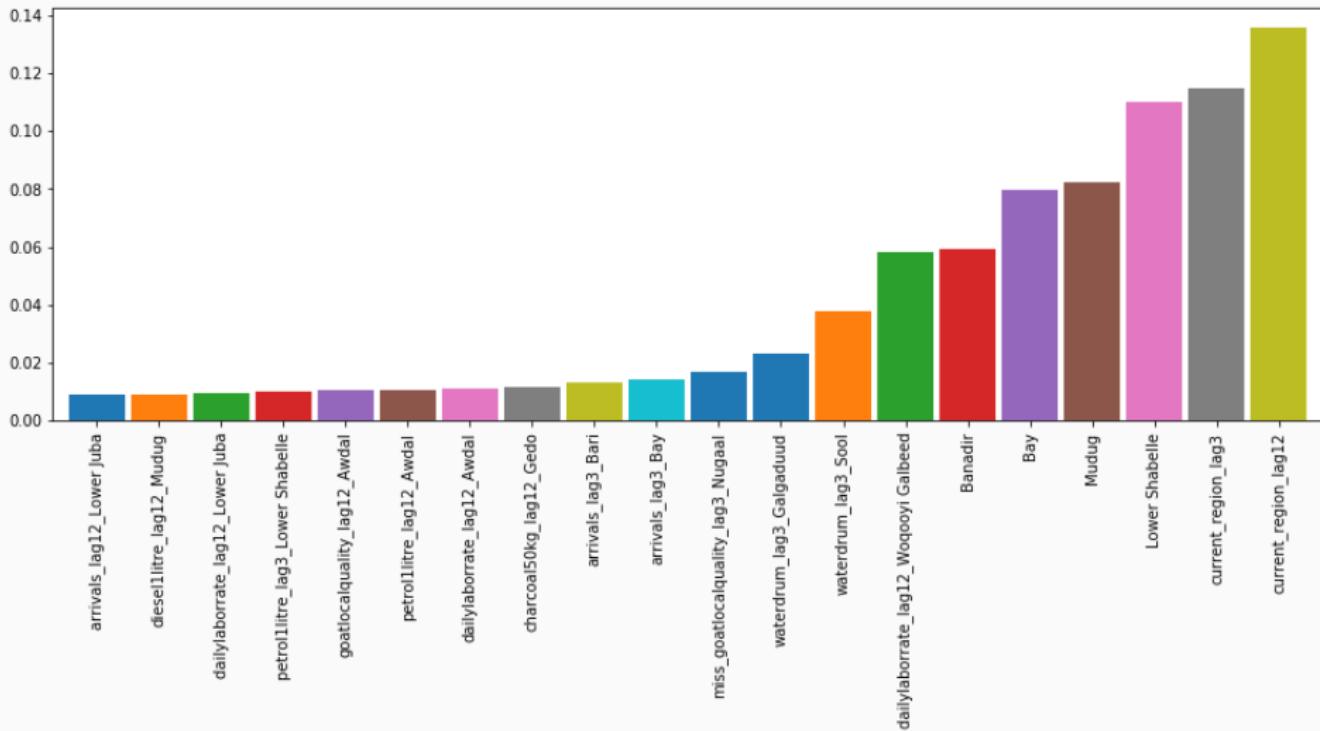
# Performance Metrics

1. Percentage of computer prediction accuracy (ACC):  $\frac{\text{predicted value}}{\text{true value}}$   
penalizes the PERCENT by which we are off
2. Mean squared error (MSE):  $(\text{true} - \text{predicted})^2$   
penalizes the AMOUNT by which we are off
3. Accuracy:  $\frac{\text{correctly classified}}{\text{all classified}}$   
penalizes the DIRECTION in which we are off

# Performance Metrics



# Explainability



# Current Priorities

## 1. Documentation

Clean, reproducible, and open-source code

## 2. Automation

Automatic data downloading and model selection

## 3. Experimentation pipeline

Tune models with clear goal in mind ;

Test different model settings and compare

## 4. Explainability

Visualization and model insights

# Questions?

**About the Data**

Open Data

The main central component of a predictive analytics project is open data. Open data as defined by [Open Data](#)

**Expand this section ↓**

**1. Introduction**  
rebecamoreno edited this page 23 hours ago - 38 revisions

**Challenge**

Project Jetson started mid-year 2017 with two operations, UNHCR Somalia and UNHCR Ethiopia Dollo Ado, Melkadida sub-office, were concerned about the ongoing drought conditions in Somalia exacerbating forced displacement. Alongside with a history of protracted conflict, Somalia experienced similar drought conditions as the year 2011, were operational emergency response was surpassed by the amount of people displaced.

Particularly, Somalia Operation reached out to UNHCR Innovation Service to be able to see if there was a way to predict forced displacement. Originally, the challenge posed by the operation was to help them building scenarios with different displacement figures, regardless of UNHCR's persons of concern (PoCs) displacement reasons. The Innovation Service explored a creative way to solve their operational challenge away from conventional methods: the use of artificial intelligence (AI) particularly, machine learning (ML) to develop predictive analytics.

**Problems related to the challenge**

Some of the problems related to the challenge to work on were:

[https://github.com/khof312/jetson\\_v1.1](https://github.com/khof312/jetson_v1.1)