# Predicting Poverty in Costa Rica
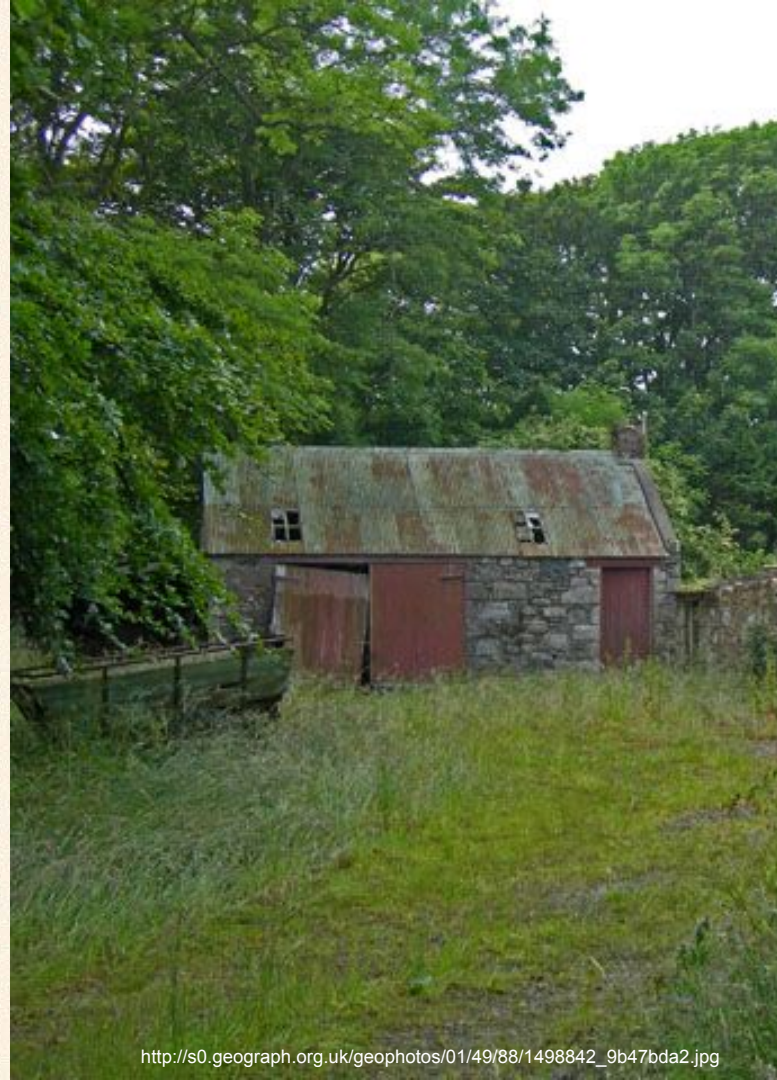
Sam Blass

Metis Classification

23 March 2022

# Objective



- Prioritize aid only to most vulnerable households
- Predict vulnerability using household observable attributes (e.g. size, dwelling quality)
- Kaggle competition 2019

http://s0.geograph.org.uk/geophotos/01/49/88/1498842_9b47bda2.jpg

# Data Overview

**Target variable** - level of poverty
- 1 = extreme poverty
- 2 = moderate poverty
- 3 = vulnerable households
- 4 = non vulnerable households

Each row corresponds to one individual

Each individual in a household gets the same household ID number

**Features** (~140)
- Individual: Age, gender, years of education
- Household: Number of household members, quality of dwelling
- Geography: Region of Costa Rica located

https://upload.wikimedia.org/wikipedia/commons/thumb/8/89/CRI_orthographic.svg/1920px-CRI_orthographic.svg.png

# Methodology

**Exploratory data analysis and feature engineering**
- Identify class imbalances
- Explore relationships not otherwise capturable in a model

**Define classification metric**
- Reflects objective of project (prioritize aid by identifying those most in need)

**Test and optimize various classification models**
- Recommend best performing model
- Train-test-split by household to avoid data leakage (Group Shuffle Split)

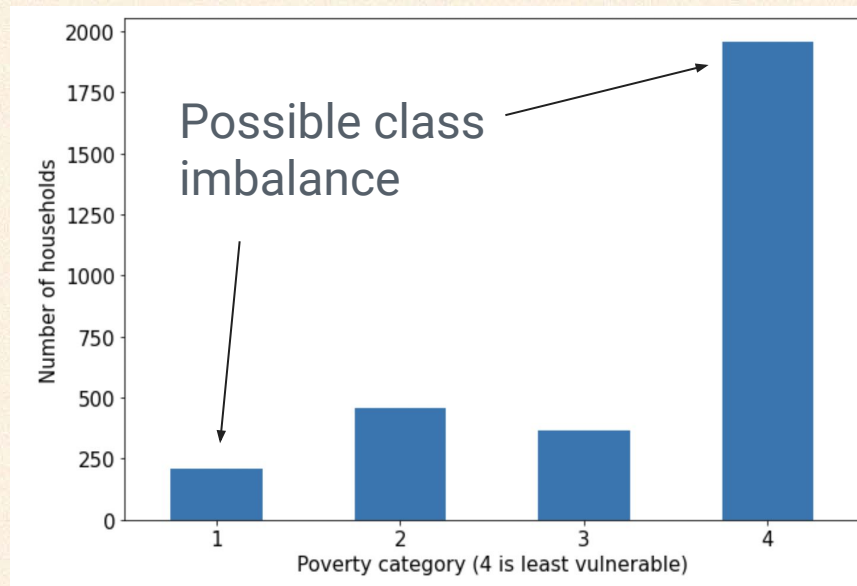# Data Prep and Model Setup

**Exploratory Data Analysis**
- Most households are not vulnerable
- Tuning model may require class balancing

**Feature Engineering**
- Scale features by geographical area to correct for cost of living

**Classification Metric: Recall**
- Minimize number of false negatives
- Prioritize category 1 (most vulnerable)
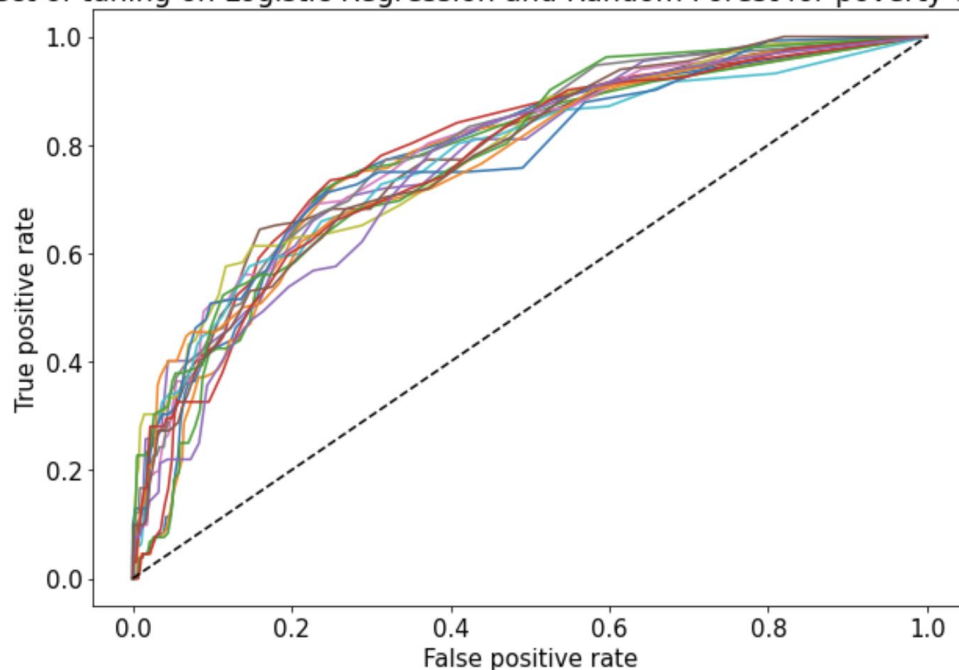- Weighted average for each recall

Possible class imbalance

# Model Testing and Tuning

| Model | Hyperparameters to tune |
|-------|------------------------|
| Random Forest | Max number of features, criterion, class weight |
| Logistic Regression | Regularization |

*Max AUC is with **unbalanced** data

| Metric | Max and Min |
|--------|-------------|
| Recall | Max: **0.345** Min: **0.273** |
| AUC | Max: **0.782** Min: **0.739** |

Effect of tuning on Logistic Regression and Random Forest for poverty category 1
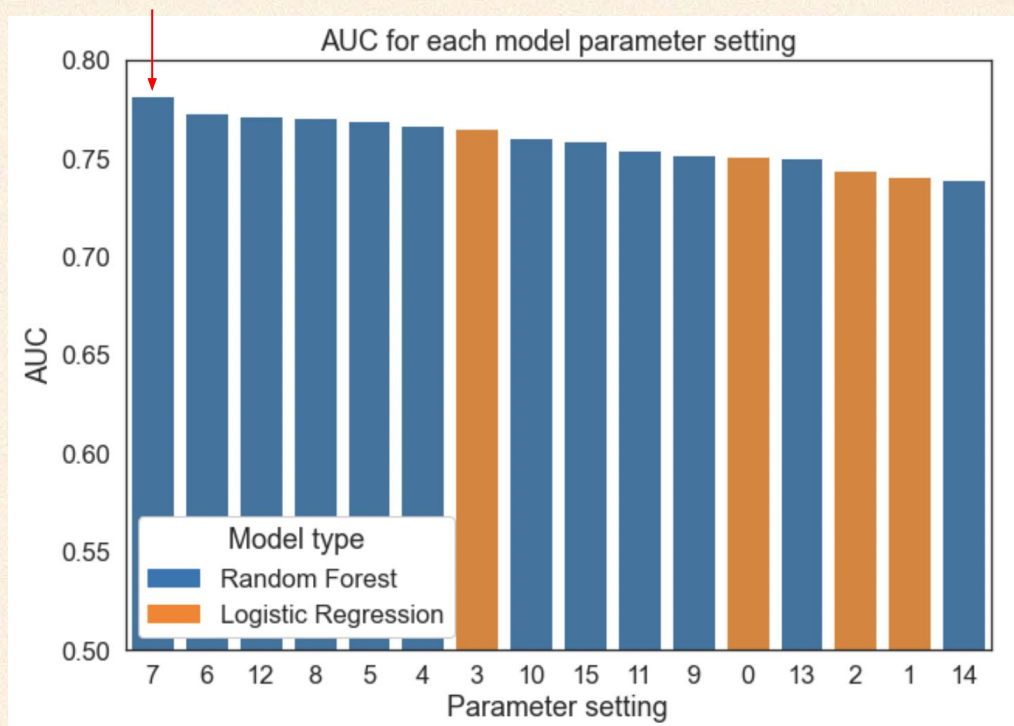
# Model Testing and Tuning

## Best results with default settings

| Model | Hyperparameters to tune |
|---|---|
| Random Forest | Max number of features, criterion, class weight |
| Logistic Regression | Regularization |

*Max AUC is with **unbalanced** data

| Metric | Max and Min |
|---|---|
| Recall | Max: **0.345** <br> Min: **0.273** |
| AUC | Max: **0.782** <br> Min: **0.739** |

Default

AUC for each model parameter setting

AUC

Model type
- Random Forest
- Logistic Regression

Parameter setting: 7 6 12 8 5 4 3 10 15 11 9 0 13 2 1 14

# Feature Importance

- Top two features are **number of females in household** and **dependency rate**
- Scaled features are important suggesting scaling relative to geographical area is important

Scaled = relative to geographic area

Top ten most important features according to Random Forest

| Feature | |
|---|---|
| Total females in household (scaled) | |
| Measure of adults to number of dependents (i.e. children, eldery) | |
| If rubbish disposal by other means (scaled) | |
| Average years of education for adults (18+) (scaled) | |
| Persons younger than 12 years of age | |
| Years of schooling | |
| If toilet connected to black hole or letrine (scaled) | |
| Owns a tablet (scaled) | |
| Age | |
| If dwelling is in precarious state (scaled) | |

Feature importance: 0.000, 0.005, 0.010, 0.015, 0.020, 0.025, 0.030, 0.035, 0.040

# Conclusions & Future Work

**01  Data Cleaning**
May improve recall scores

**02  Hyperparameter tuning**
No improvement in recall

**03  Feature importance**
Identified which features to examine further

**04  Target variable subjectivity**
Likely reduces model performance

# Thanks!