# Machine Learning Methods in Vulnerable Communities Classification

- Carlos Gomes de Oliveira
  DBA and Data Engineer

  Linkedin.com/in/CarlosOliveira

# Applying Machine Learning Methods in People Classification

Carlos Gomes de Oliveira*, Eduardo Terra Morelli†, and Raul Senna Ferre
Instituto INFNET
Escola Superior da Tecnologia da Informação
Rio de Janeiro, RJ
Email: *czgrqg@gmail.com, †emorelli1966@gmail.com, ‡raul.ferreira@prof.infnet.e

*Abstract*—**The missing people database offers several challenges regarding its analysis. The database may be incomplete since usually there is not a commitment from the authorities to keep consolidated and detailed information on missing persons. There are many apparent uncorrelated attributes, as different persons have different attribute groups. Many of these records contain attributes with incorrect and missing information.**

**This paper intends to demonstrate the use of machine learning methods as Decision Trees, Random Forests, Logistic Regression**
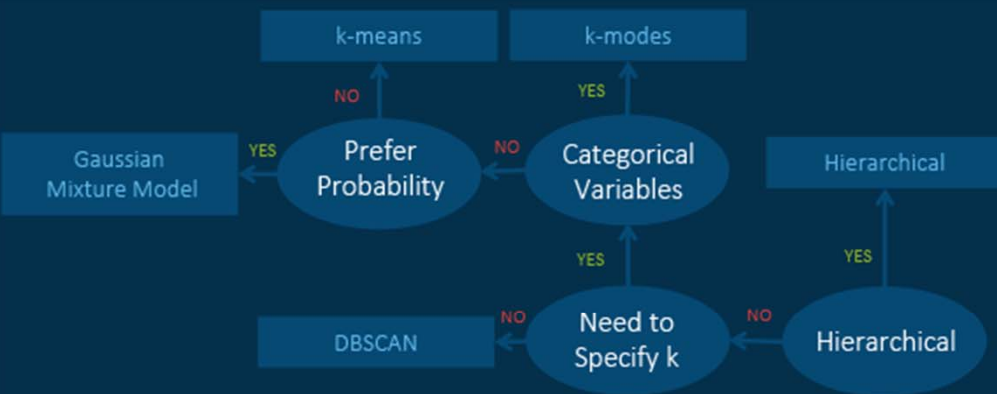
community represented by the data, th
tics that distinguish them. This can be
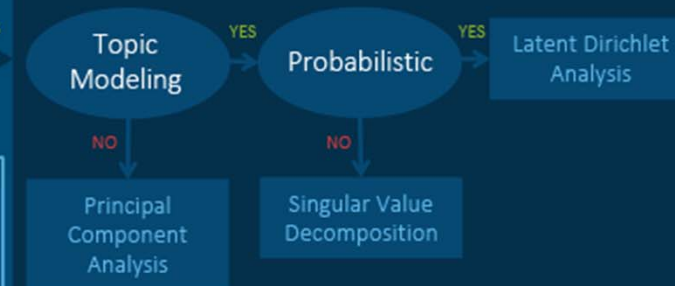machine learning methods.

## II. RELATED WO

Machine learning algorithms are wid
for example, Cornel and Mirela us
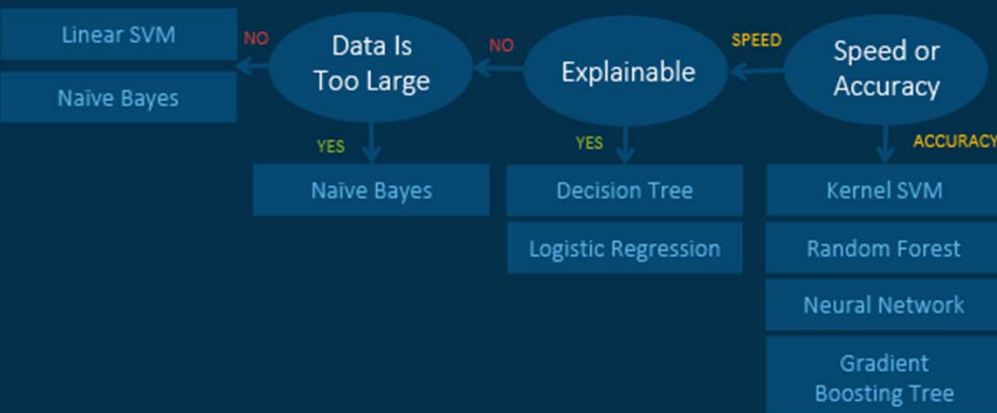
# Machine Learning Algorithms Cheat Sheet
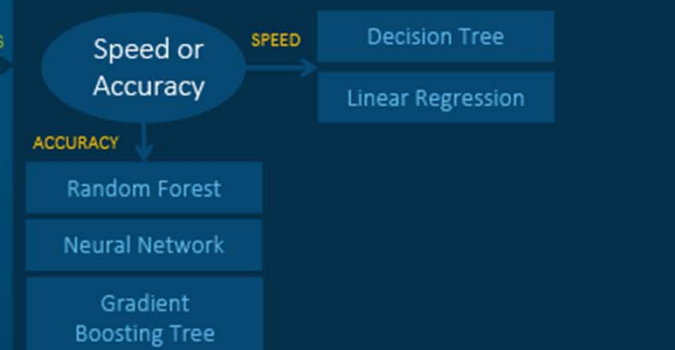
## Unsupervised Learning: Clustering

k-means | k-modes

Gaussian Mixture Model ← YES — Prefer Probability ← NO — Categorical Variables

NO (k-means) / YES (k-modes)

Hierarchical

DBSCAN ← NO — Need to Specify k ← NO — Hierarchical

YES / YES

## START

Dimension Reduction — YES → Topic Modeling — YES → Probabilistic — YES → Latent Dirichlet Analysis

NO

Have Reponses

NO / YES

## Unsupervised Learning: Dimension Reduction

Topic Modeling — NO → Principal Component Analysis

Probabilistic — NO → Singular Value Decomposition

## Supervised Learning: Classification

Linear SVM
Naïve Bayes

Data Is Too Large ← NO — Explainable ← NO — Speed or Accuracy ← Predicting Numeric

SPEED

YES → Naïve Bayes

YES → Decision Tree / Logistic Regression

ACCURACY → Kernel SVM / Random Forest / Neural Network / Gradient Boosting Tree

## Supervised Learning: Regression

Predicting Numeric — YES → Speed or Accuracy

SPEED → Decision Tree / Linear Regression

ACCURACY → Random Forest / Neural Network / Gradient Boosting Tree

# scikit-learn algorithm cheat-sheet

**START**

## classification

- kernel approximation
- SVC Ensemble Classifiers
- KNeighbors Classifier
- SGD Classifier
- Naive Bayes
- Text Data
- Linear SVC
- <100K samples

- **>50 samples** — NO → get more data
- YES → predicting a category
- predicting a category — YES → do you have labeled data
- do you have labeled data — YES → <100K samples
- <100K samples — YES → Linear SVC
- Text Data — YES → Naive Bayes; NO → KNeighbors Classifier
- NOT WORKING → SGD Classifier → NOT WORKING → kernel approximation
- NOT WORKING → SVC Ensemble Classifiers

## regression

- SGD Regressor
- ElasticNet Lasso
- SVR(kernel='rbf') EnsembleRegressors
- <100K samples
- few features should be important
- RidgeRegression SVR (kernel='linear')

- predicting a **quantity** — YES → <100K samples
- <100K samples — NO → SGD Regressor; YES → few features should be important
- few features should be important — YES → ElasticNet Lasso; NO → RidgeRegression SVR (kernel='linear')
- NOT WORKING → SVR(kernel='rbf') EnsembleRegressors

## clustering

- Spectral Clustering GMM
- KMeans
- <10K samples
- MiniBatch KMeans
- number of categories known
- <10K samples
- MeanShift VBGMM

- labeled data — NO → number of categories known
- number of categories known — YES → <10K samples; NO → <10K samples
- <10K samples — YES → KMeans; NO → MiniBatch KMeans
- KMeans — NOT WORKING → Spectral Clustering GMM
- <10K samples — YES → MeanShift VBGMM; NO → tough luck

## dimensionality reduction

- Randomized PCA
- Isomap Spectral Embedding
- LLE
- <10K samples
- kernel approximation

- predicting a category — NO → predicting a quantity
- predicting a quantity — NO → just looking
- just looking — YES → Randomized PCA; NO → predicting structure
- Randomized PCA — NOT WORKING → <10K samples
- <10K samples — YES → Isomap Spectral Embedding; NO → kernel approximation
- Isomap Spectral Embedding — NOT WORKING → LLE
- predicting structure → tough luck

# Objective

- Identify risk communities in a community using machine learning methods
    - Decision Trees
    - Random Forests
    - XGBoost
    - Logistic Regression

# Machine Learning

Basic equation: $Obj(\theta) = L(\theta) + \Omega(\theta)$
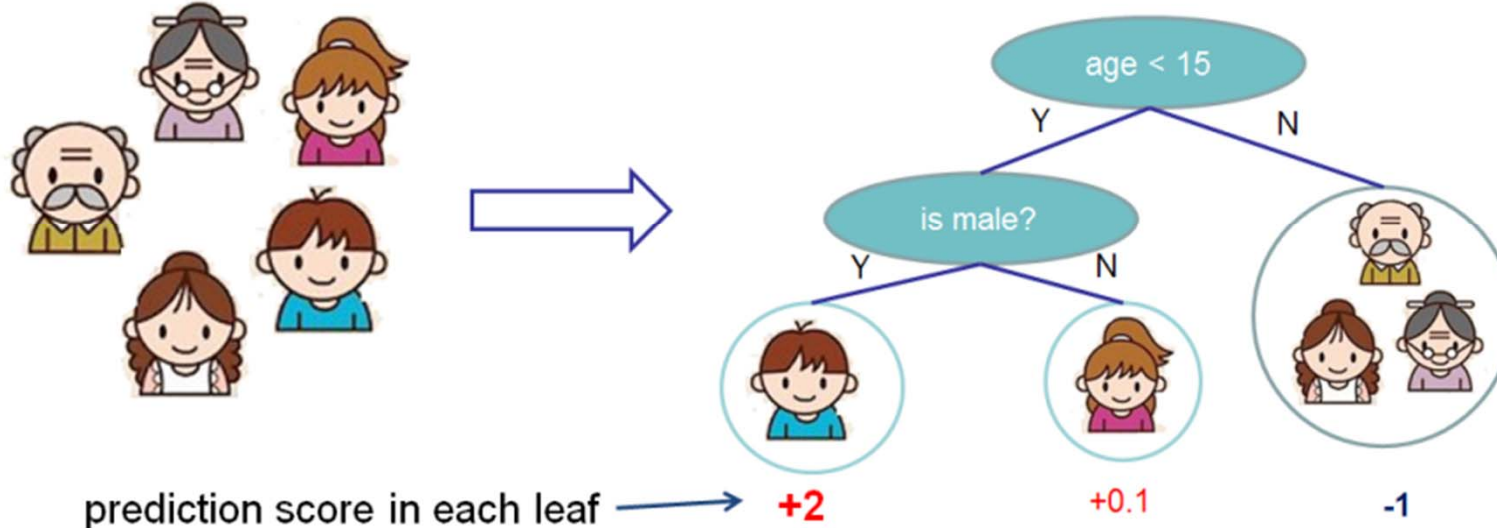
- L – Training Loss

- $\Omega$ – Regularization

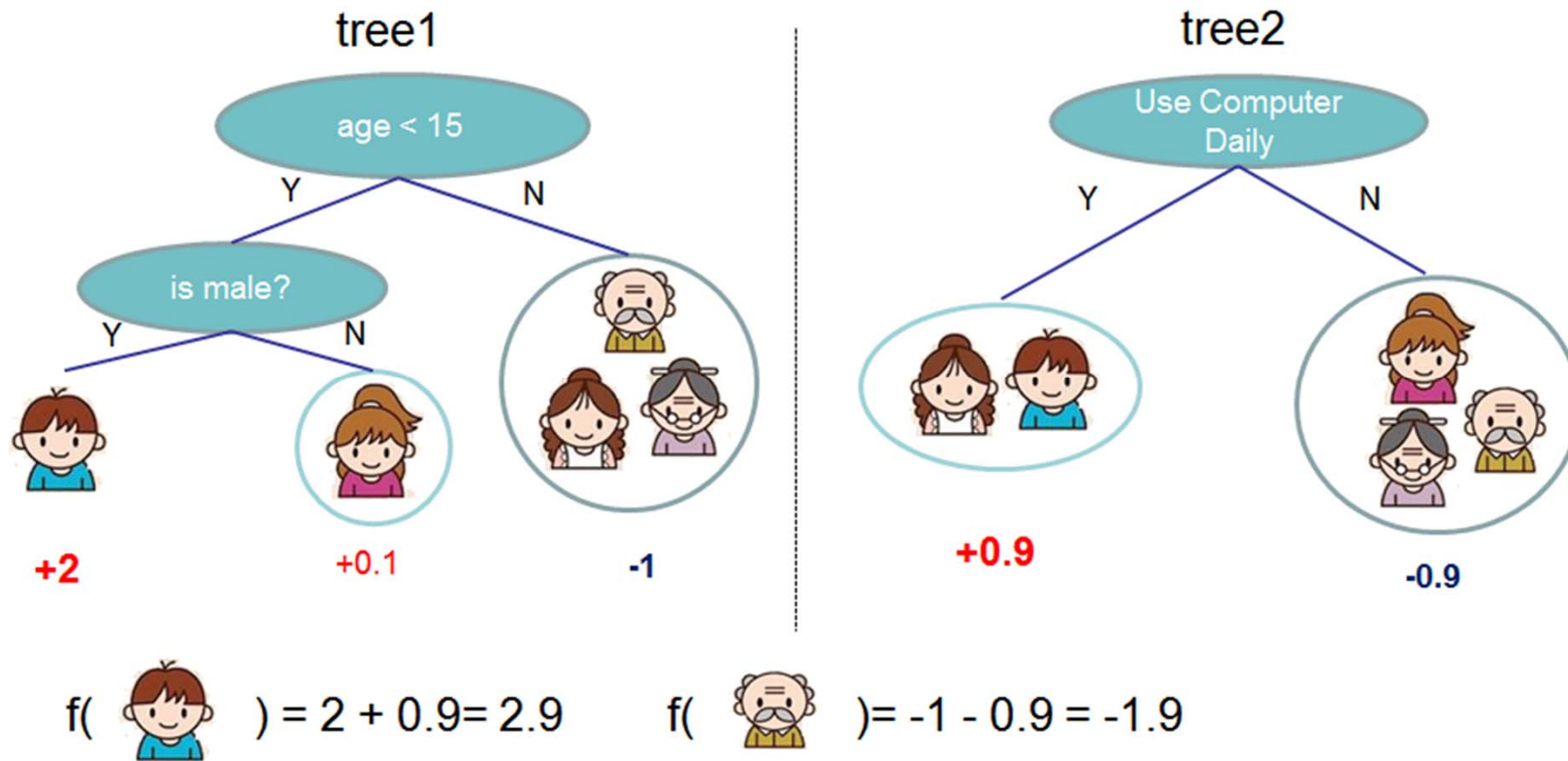    Ex: Linear regression, Logistic Regression, Decision Trees

# Decision Trees



Input: age, gender, occupation, ...

Does the person like computer games

prediction score in each leaf → **+2**   +0.1   -1

# Random Forests

# Gradient Boosting (XGBoost)

Instance index    gradient statistics

1    g1, h1

2    g2, h2

3    g3, h3

4    g4, h4

5    g5, h5

age < 15

Y    N

is male?

Y    N

$I_3 = \{2, 3, 5\}$
$G_3 = g_2 + g_3 + g_5$
$H_3 = h_2 + h_3 + h_5$

$I_1 = \{1\}$
$G_1 = g_1$
$H_1 = h_1$

$I_2 = \{4\}$
$G_2 = g_4$
$H_4 = h_4$

$$Obj = -\sum_j \frac{G_j^2}{H_j + \lambda} + 3\gamma$$
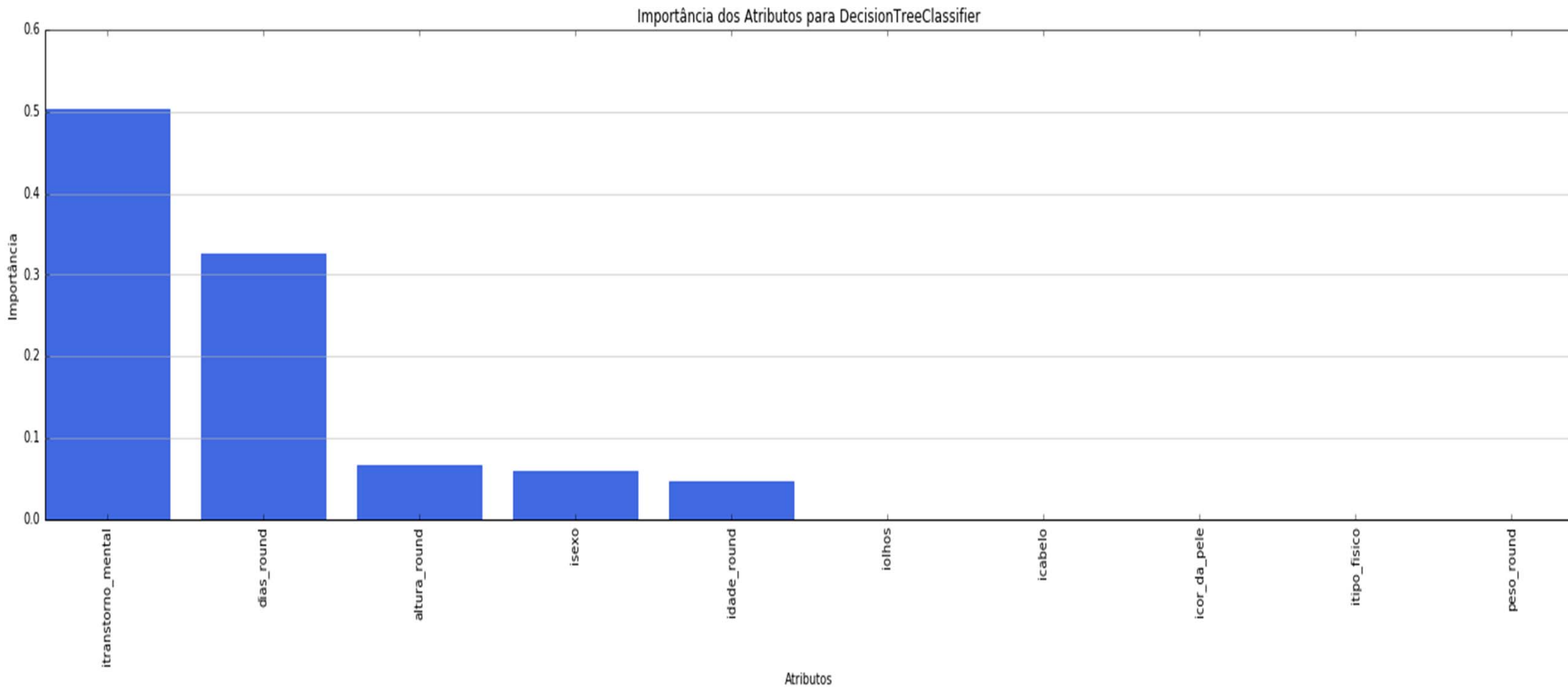
The smaller the score is, the better the structure is

g1, h1    g4, h4      g2, h2    g5, h5    g3, h3

$G_L = g_1 + g_4$      $G_R = g_2 + g_3 + g_5$

## FEATURES USED IN MODEL

| Feature Name | Feature Type |
|---|---|
| Days Missing | Quantitative |
| Height | Quantitative |
| Weight | Quantitative |
| Age | Quantitative |
| Physical Type | Categorical |
| Skin Color | Categorical |
| Eye Color | Categorical |
| Sex | Categorical |
| Hair Color | Categorical |
| Mental Impairment | Categorical |

## MODEL ACCURACY

| Algorithm | Test#1 | Test#2 | Test#3 | Test#4 |
|---|---|---|---|---|
| Decision Trees | 69.82% | 68.39% | 68.68% | 67.24% |
| Random Forest | **80.46%** | 79.02% | **81.32%** | 77.58% |
| Logistic Regression | 75.86% | 72.70% | 75.29% | 71.55% |
| XGBoost | **81.03%** | 79.02% | 79.31% | 77.59% |

# Decision Tree Classifier 69%



Importância dos Atributos para DecisionTreeClassifier

# LogisticRegression 73%



Importância dos Atributos para LogisticRegression

# RandomForestClassifier 77%



Importância dos Atributos para RandomForestClassifier
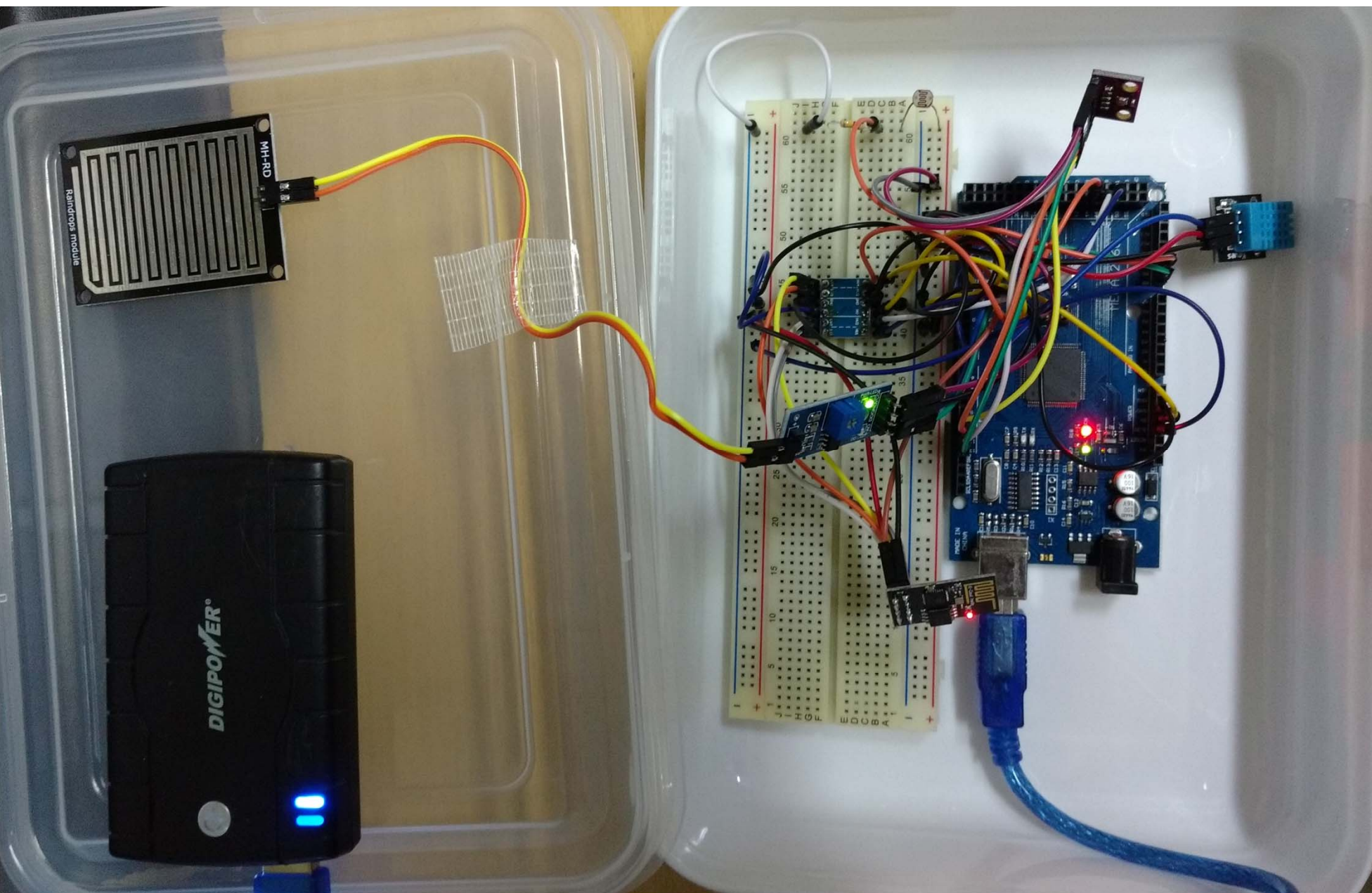
# XGBClassifier 79%



Importância dos Atributos para XGBClassifier

# OTHER
# PROJECTS

# QUESTIONS?