emata

# Don't Get Lost in the Random Forests: A Beginner's Guide

Moses Bomera, Emata

April 29th, 2021

# Moses Bomera
## Data Scientist
## Emata (& Laboremus)

Works on alternative credit scoring for microfinance institution Emata, a spin-off of fintech company Laboremus.

Prior to joining Laboremus, worked on research in computer vision and natural language processing at netLabs!UG Research Centre, Makerere University.

**We are on a mission to provide affordable digital loans to millions of farmers in East Africa**

Today, agri-financing in Africa does NOT work.

**Farmers are risky clients**

NO (CREDIT) HISTORY, NO COLLATERAL

**Farmers need small loans**
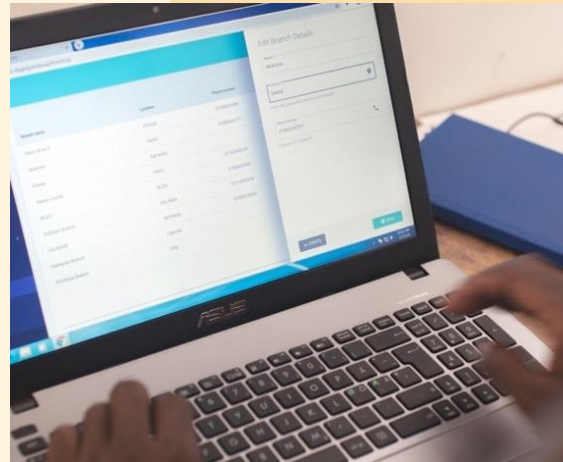
HIGH OPERATING COSTS
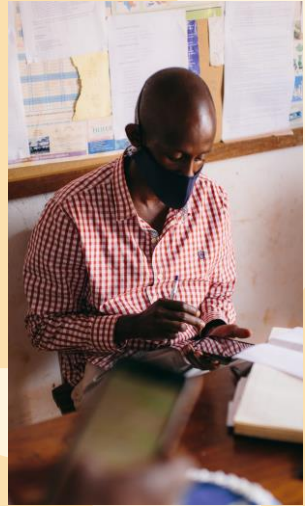
**50-100%**

ANNUAL INTEREST RATE

# What does **Emata** do?

1. Digitise agricultural cooperatives to get data

2. Turn data into credit scoring
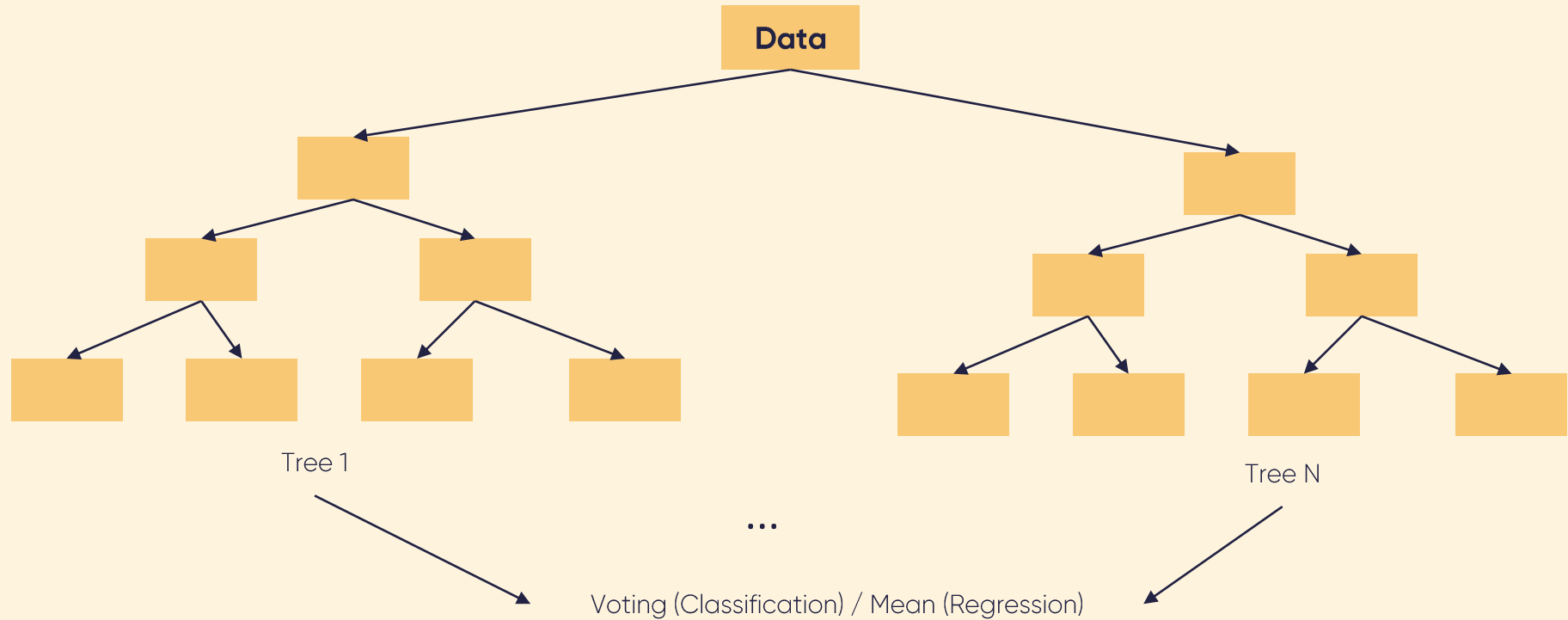
3. Offer digital and affordable loans to farmers

# Agenda

emata

# What is a Random Forest?

# What is a Random Forest?

It is an ensemble of decision trees.



Data

Tree 1

Tree N

...

Voting (Classification) / Mean (Regression)

# Random Forests

## A decision tree (1/3)

Did you sign up for DataFest?

True

False

Attend the Random Forest Talk.

Can't attend the Random Forest Talk

# Random Forests

## A decision tree (2/3)

```
        ┌─────────────────┐
        │      Body       │
        │ temperature >   │
        │      38°C       │
        └─────────────────┘
         True         False
         ↓               ↓
┌─────────────────┐  ┌─────────────────┐
│ Huh! Corona??*  │  │  No Worries.*   │
└─────────────────┘  └─────────────────┘
```

*For tutorial purpose only.

# Random Forests

## A decision tree (3/3)

# Dataset

| petal length | petal width | target |
|---|---|---|
| 1.4 | 0.2 | setosa |
| 1.3 | 0.2 | setosa |
| ... | ... | ... |
| 5.1 | 1.8 | virginica |

[Iris dataset](#)

Left; Iris setosa,
Top;  Iris virginica,
Bottom; Iris veriscolor

**The Theory**

petal length (cm) <= 2.45
gini = 0.667
samples = 150
value = [50, 50, 50]
class = setosa

True

False

gini = 0.0
samples = 50
value = [50, 0, 0]
class = setosa

petal width (cm) <= 1.75
gini = 0.5
samples = 100
value = [0, 50, 50]
class = versicolor

gini = ??
samples = 54
value = [0, 49, 5]
class = versicolor

gini = 0.043
samples = 46
value = [0, 1, 45]
class = virginica

# Node Features

**Samples:** the number of training instances that particular node applies to.

**Gini:** measures how pure (or impure) a node is.

**Class:** the target represented by a given node.

$p_{i,k}$ is the ratio of class k instances to all the training instances in $i^{th}$ node.

gini = 0.043
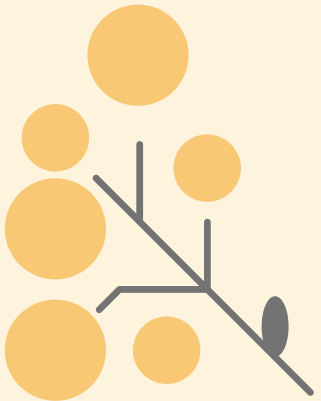samples = 46
value = [0, 1, 45]
class = virginica

$$G_i = 1 - \sum_{k=1}^{n} p_{i,k}^{2}$$

# Calculating the gini value

gini = ??
samples = 54
value = [0, 49, 5]
class = versicolor

$$G_i = 1 - \sum_{k=1}^{n} p_{i,k}^2$$

$$1 - \left(\frac{0}{54}\right)^2 - \left(\frac{49}{54}\right)^2 - \left(\frac{5}{54}\right)^2 =$$
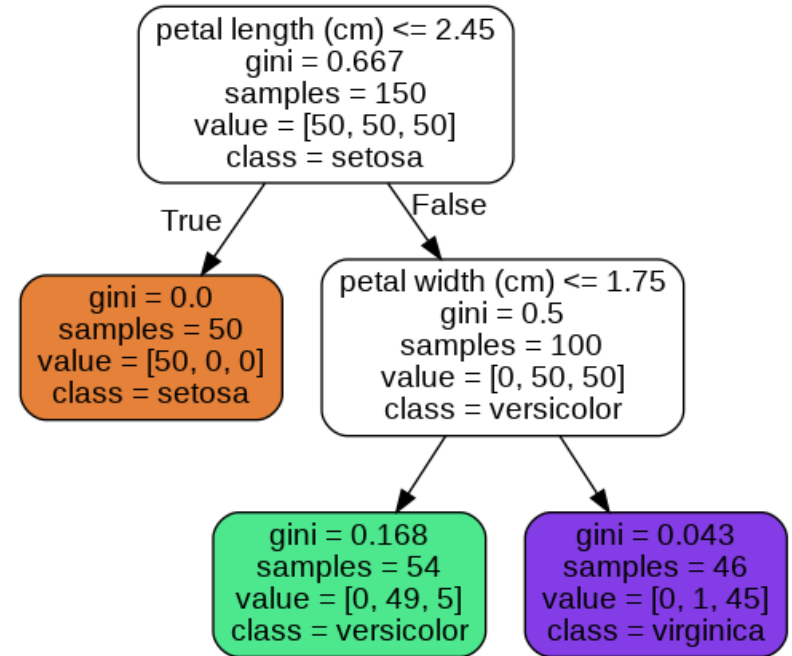
# Classification and Regression Tree (CART) Algorithm

$$J(k, t_k) = \frac{m_{\text{left}}}{m} G_{\text{left}} + \frac{m_{\text{right}}}{m} G_{\text{right}}$$
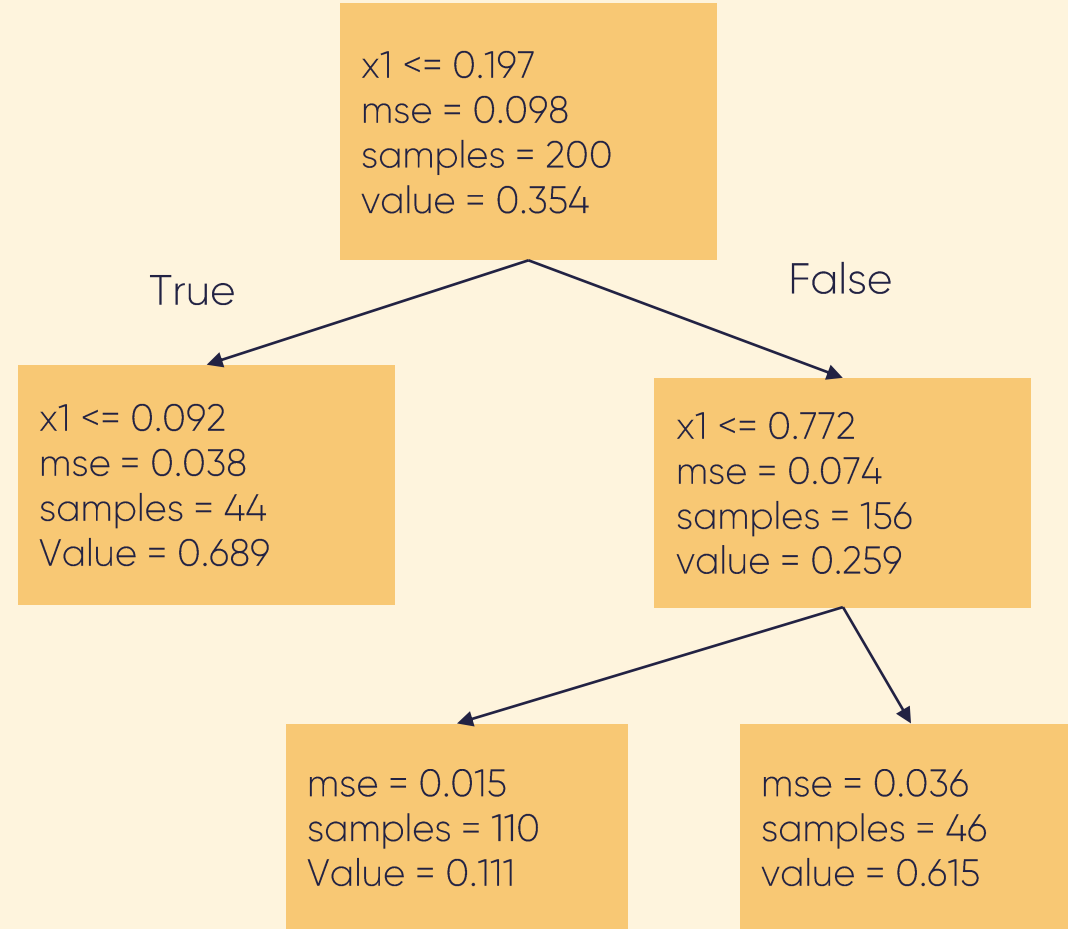
where $\begin{cases} G_{\text{left/right}} & \text{measures the impurity of the left/right subset,} \\ m_{\text{left/right}} & \text{is the number of instances in the left/right subset.} \end{cases}$

1. CART splits the training set into two subsets using a single feature, k and a threshold purity, $t_k$.
2. The selection process involves minimizing the cost function
3. 1 and 2 are repeated recursively for each subset until the maximum depth is reached.

petal length (cm) <= 2.45
gini = 0.667
samples = 150
value = [50, 50, 50]
class = setosa

True

False

gini = 0.0
samples = 50
value = [50, 0, 0]
class = setosa

petal width (cm) <= 1.75
gini = 0.5
samples = 100
value = [0, 50, 50]
class = versicolor

gini = 0.168
samples = 54
value = [0, 49, 5]
class = versicolor

gini = 0.043
samples = 46
value = [0, 1, 45]
class = virginica

# Regression

- Similar to the classification approach.

- The prediction is the average of the samples associated with the leaf node.

- For the gini impurity split, the regression tree uses mean squared error i.e. attempts to minimize the mse.

x1 <= 0.197
mse = 0.098
samples = 200
value = 0.354

True

False

x1 <= 0.092
mse = 0.038
samples = 44
Value = 0.689

x1 <= 0.772
mse = 0.074
samples = 156
value = 0.259

mse = 0.015
samples = 110
Value = 0.111

mse = 0.036
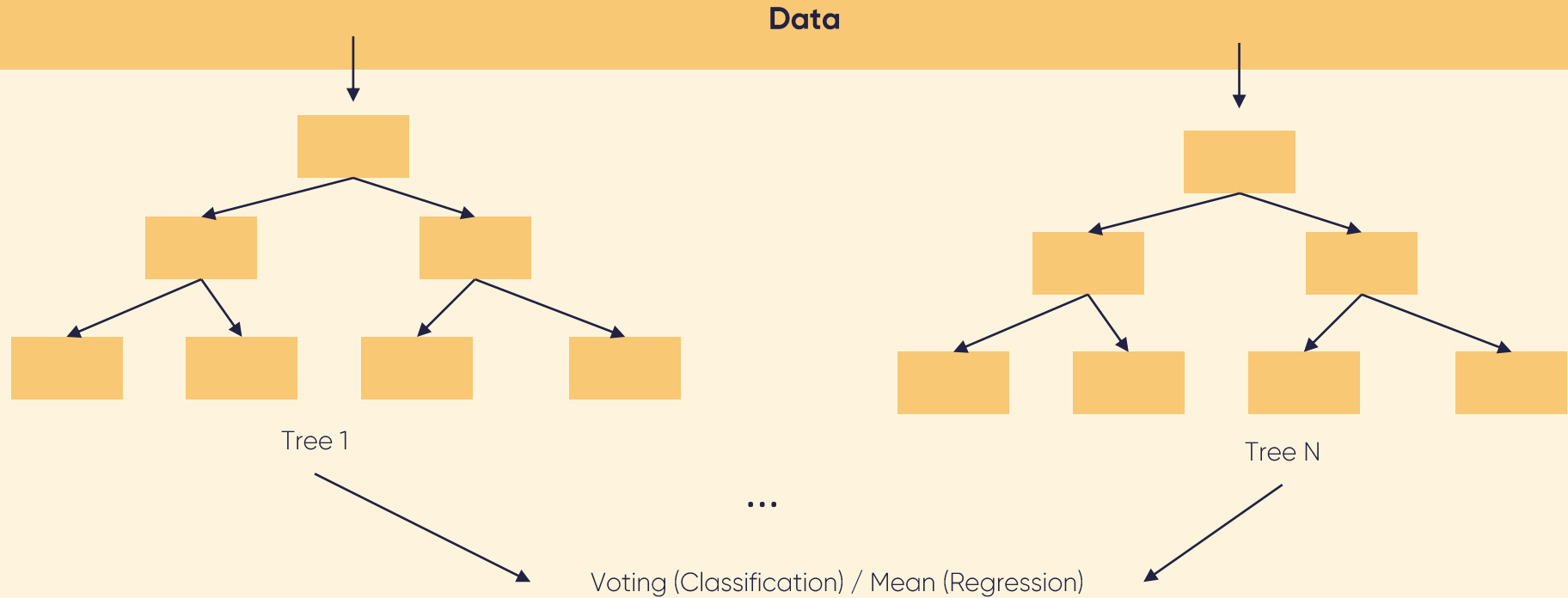samples = 46
value = 0.615

# Gini impurity or Entropy?

- **While gini is the go-to impurity measure**, you can also use entropy as a measurement.

- Entropy is based on the same concept from thermodynamics.

- Entropy approaches zero when molecules are still and well ordered.

- In machine learning, entropy is zero when a set contains instances of only one class.

$$H_i = - \sum_{\substack{k=1 \\ p_{i,k} \neq 0}}^{n} p_{i,k} \log_2 \left( p_{i,k} \right)$$

# What is a Random Forest?

It is an ensemble of decision trees.

**Data**

Tree 1

Tree N

...
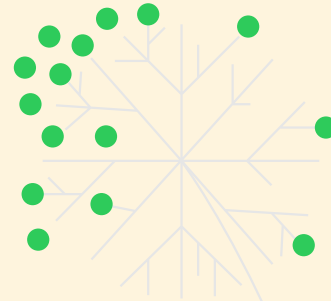
Voting (Classification) / Mean (Regression)

# How does it work? (1/2)

- A random forest is an ensemble of decision trees.

- Trained usually through the bagging method (or pasting).

- Each tree in a random forest gives a prediction, for classification, the class with the majority of votes is the prediction.

- In regression, the average of each tree's prediction is the prediction.

# How does it work? (2/2)

- RF introduces randomness when growing trees.

- Doesn't make a split using the very best feature, instead it selects the best feature among a random subset of features.

- Results in greater tree diversity.

- Trading a higher bias with lower variance yielding a better model than an individual decision tree.

# Regularization Hyperparameters

- *n_estimators* – the number of trees to use.

- *max_depth* – the depth of the tree.

- *min_samples_split* – the minimum number of samples an internal node must have before it can be split.

- *min_samples_leaf* – the minimum number of samples required to be at a leaf node

- *max_leaf_nodes* – the maximum number of leaf nodes.

- *max_features* – the maximum number of features that are evaluated for splitting at each node.

Increasing *min_\** hyperparameters or reducing *max_\** hyperparameters will regularize the model.
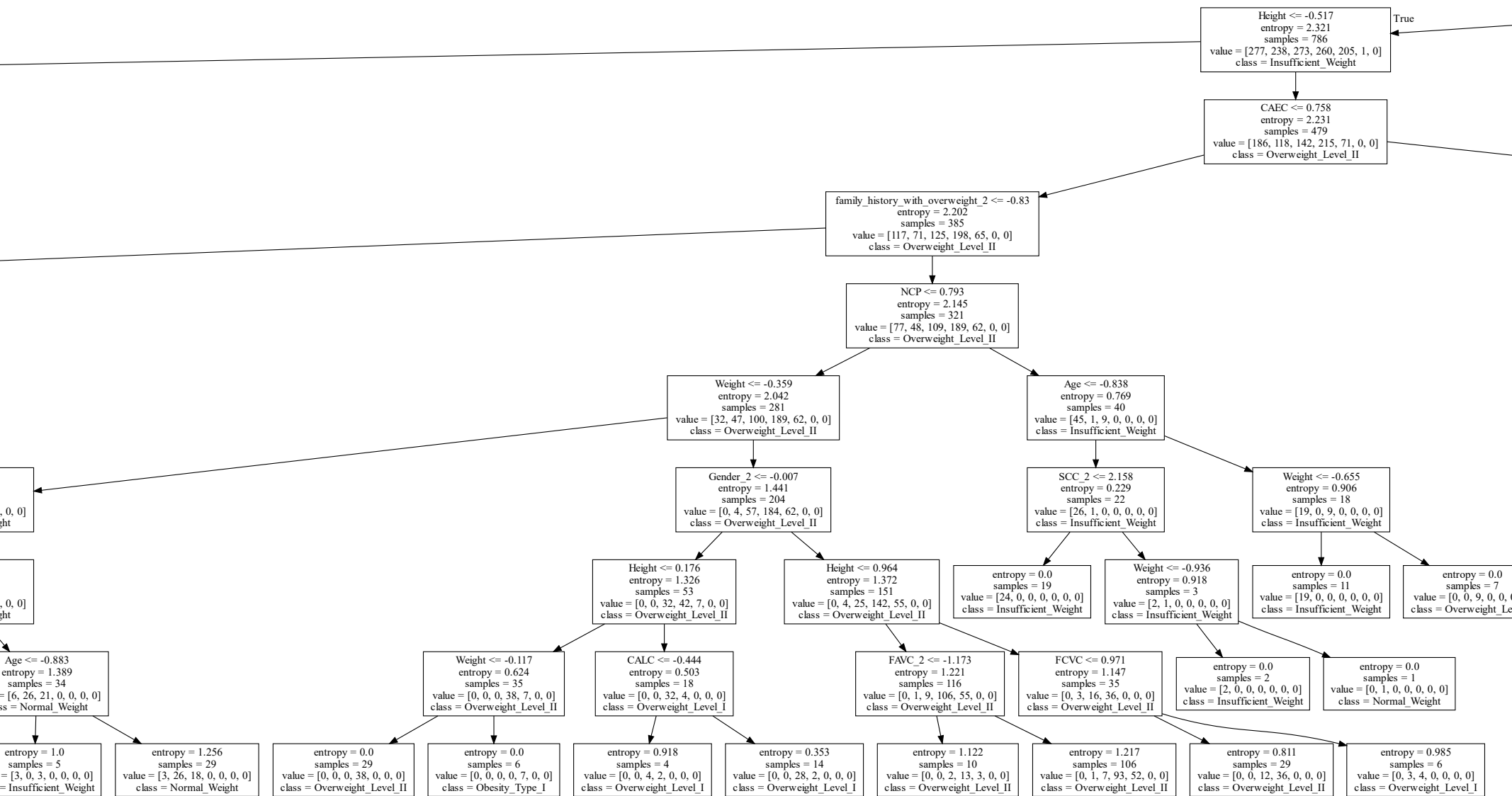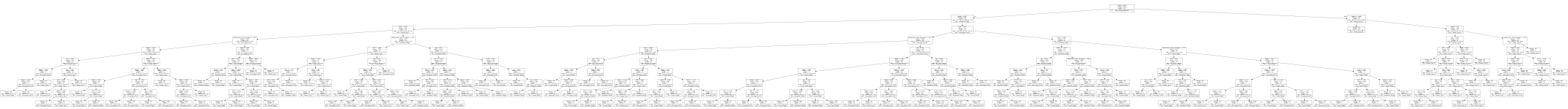
# TLDR;

**PROS of using Random Forests**

- Random Forests can be used for both classification and regression.

- Can handle large datasets with high dimensionality, (can be used for dimensionality reduction).

- Works well with unscaled datasets.

**CONS of using Random Forests**

- It is a high variance model, so it is important to regularize the model when training.

- Given the numerous hyper-parameters, and the number of trees, the decision process becomes a bit of a black box.

# Tutorial

[Tutorial Notebook](Tutorial Notebook)

# Reference Material

1. Chapter 6, 7 Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow by Aurélien Géron

2. https://ieeexplore.ieee.org/abstract/document/598994 – Random Forest academic paper.

3. https://towardsdatascience.com/understanding-random-forest-58381e0602d2

4. https://github.com/kjw0612/awesome-random-forest – Projects that have used Random Forests to achieve amazing solutions.

5. https://www.quora.com/What-are-the-advantages-and-disadvantages-for-a-random-forest-algorithm – Discussions on the advantages and disadvantages of Random Forests

6. https://www.youtube.com/watch?v=7VeUPuFGJHk – StatQuest: Decision Trees

7. https://www.youtube.com/watch?v=J4Wdy0Wc_xQ – StatQuest: Random Forests Part 1- Building, Using and Evaluating

# Thank you!

Questions?

# emata

The future of farmer financing
www.emata.ug