**University of Minho**
School of Engineering

# Machine Learning and Decision-Making

ADI @ LEI/3º, MiEI/4º - 2º Semestre

Filipe Gonçalves, Inês Alves, Cesar Analide

Part VI – March 2022

# Contents

| 2 | Tree-based models | Loops | Hands On |
|---|---|---|---|

- The Learner-Predictor concept and Tree-based models

- (Loops)

- (Model Tuning)

- Hands On

# A ML Model Goes to a Job Interview

| Tree-based models | Loops | Hands On |
|---|---|---|

*Interviewer*: What's your biggest strength?
*ML Candidate*: I learn very well!

# A ML Model Goes to a Job Interview

| Tree-based models | Loops | Hands On

*Interviewer*: What's your biggest strength?
*ML Candidate*: I learn very well!

*Interviewer*: Ok! So, what's 20+15?
*ML Candidate* : It's 5.

# A ML Model Goes to a Job Interview

| Tree-based models | Loops | Hands On |
|---|---|---|

*Interviewer*: What's your biggest strength?
*ML Candidate*: I learn very well!

*Interviewer*: Ok! So, what's 20+15?
*ML Candidate* : It's 5.

*Interviewer*: Not even close. It's 35.
*ML Candidate* : It's 20.

# A ML Model Goes to a Job Interview

| Tree-based models | Loops | Hands On |
|---|---|---|

*Interviewer*: What's your biggest strength?
*ML Candidate*: I learn very well!

*Interviewer*: Ok! So, what's 20+15?
*ML Candidate* : It's 5.

*Interviewer*: Not even close. It's 35.
*ML Candidate* : It's 20.
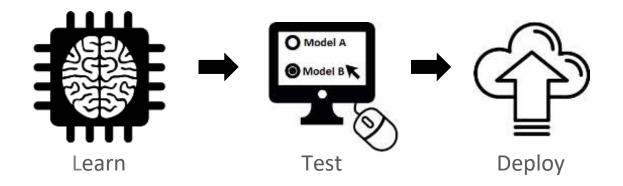
*Interviewer*: I said 35.
*ML Candidate* : 33.

# A ML Model Goes to a Job Interview

| Tree-based models | Loops | Hands On |
|---|---|---|

*Interviewer*: What's your biggest strength?
*ML Candidate*: I learn very well!

*Interviewer*: Ok! So, what's 20+15?
*ML Candidate* : It's 5.

*Interviewer*: Not even close. It's 35.
*ML Candidate* : It's 20.

*Interviewer*: I said 35.
*ML Candidate* : 33.

*Interviewer*: It's 35.
*ML Candidate* : It's 35.

*Interviewer*: Hired!

# The Learner-Predictor Concept

| **TREE-BASED MODELS** | Loops | Hands On |
|---|---|---|

Supervised algorithms imply a learning phase before applying the model to new data. But they also require a testing phase and a tuning phase!

In KNIME we implement supervised algorithms with Learner and Predictor nodes

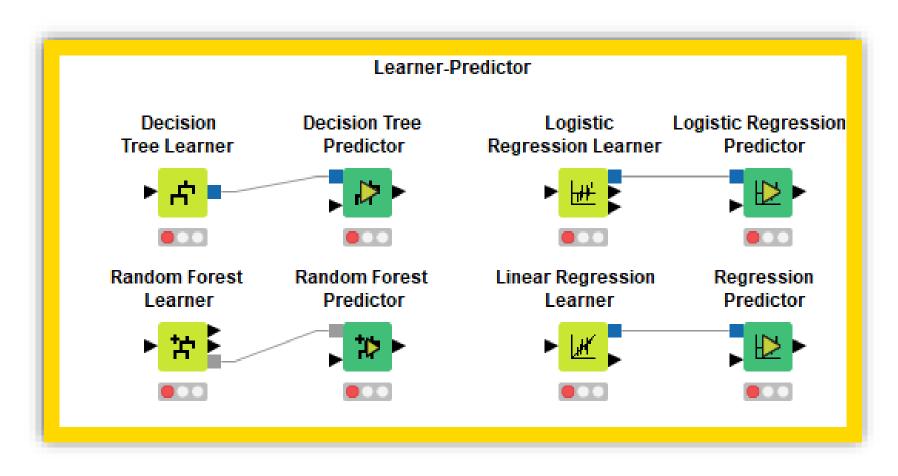Learn       Test       Deploy

# Learner-Predictor

# Some Learner-Predictor Nodes

# Decision Trees

**TREE-BASED MODELS**          Loops          Hands On



ROOT Node

Branch/ Sub-Tree

**Splitting**

Decision Node

A Decision Node

Terminal Node

Decision Node

Terminal Node

Terminal Node

B

C

Terminal Node

Terminal Node

**Note:-** A is parent node of B and C.

# Decision Trees

# Tuning Numeric Parameters

# Tuning Numeric Parameters
## Parameter Optimization Loop Nodes

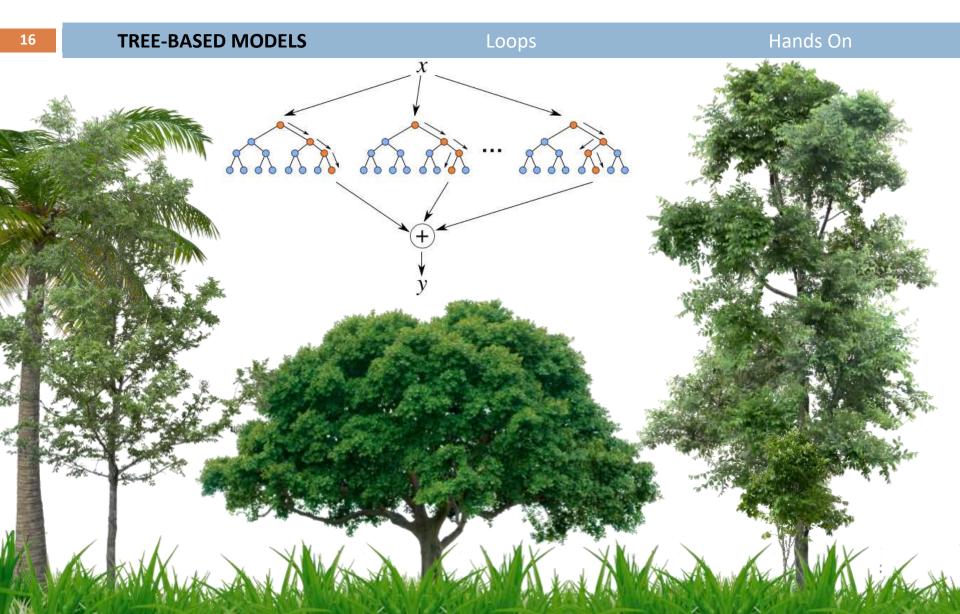Optimize the value of some parameters with respect to a cost function

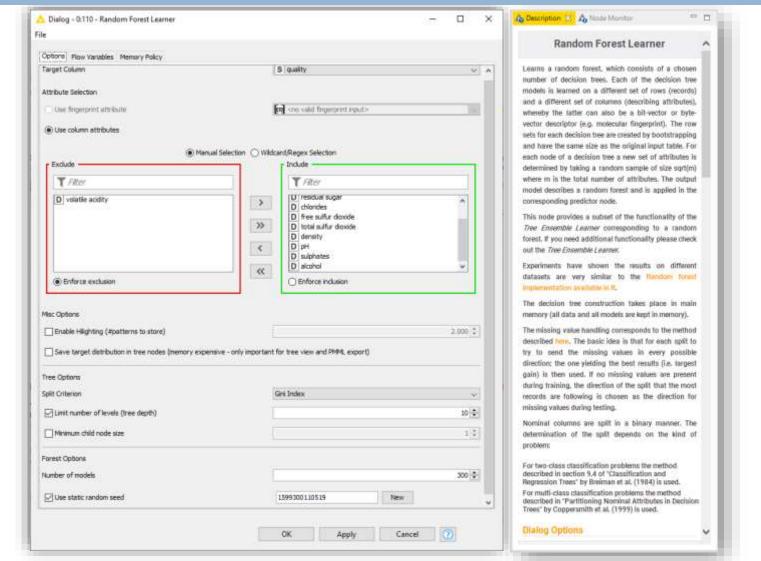# Model Hyper-parameters

Revelant concepts…

- Model Parameters: a model's (internal) configuration variable whose value is estimated from training data, i.e., the value is not set manually. Some examples include:
    - Weights in Artificial Neural Networks
    - Support vectors in Support Vector Machines

- Model Hyperparameters: a model's (external) configuration variable whose value can be set manually. It is difficult to know, beforehand, the best value of each hyperparameter. Tuning a model consists in finding the best (or, at least, a good) configuration of hyperparameters. Examples include:
    - Optimizer and learning rate in Artificial Neural Networks
    - C and sigma in Support Vector Machines
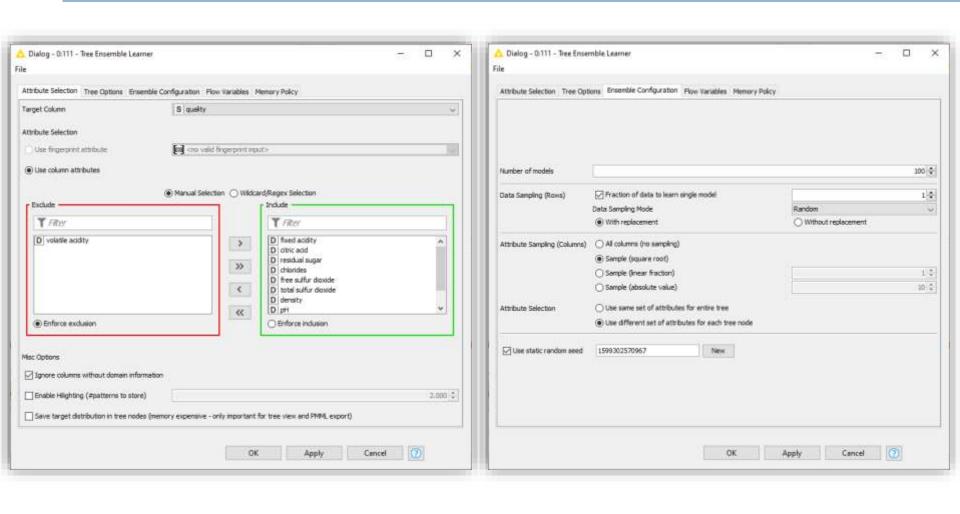    - Quality measure and pruning method in Decision Trees

# Random Forests

**TREE-BASED MODELS**          Loops          Hands On

# Random Forests

# Tree Ensembles

# Hands On