



MACHINE LEARNING WITH TREE-BASED MODELS IN R

**Welcome to the  
course!**

Erin LeDell & Gabriela de Queiroz

Machine Learning Scientist & Data Scientist



# Tree-based models

- Interpretability + Ease-of-Use + Accuracy
- Make Decisions + Numeric Predictions



# What you'll learn:

- Interpret and explain decisions
- Explore different use cases
- Build and evaluate classification and regression models
- Tune model parameters for optimal performance

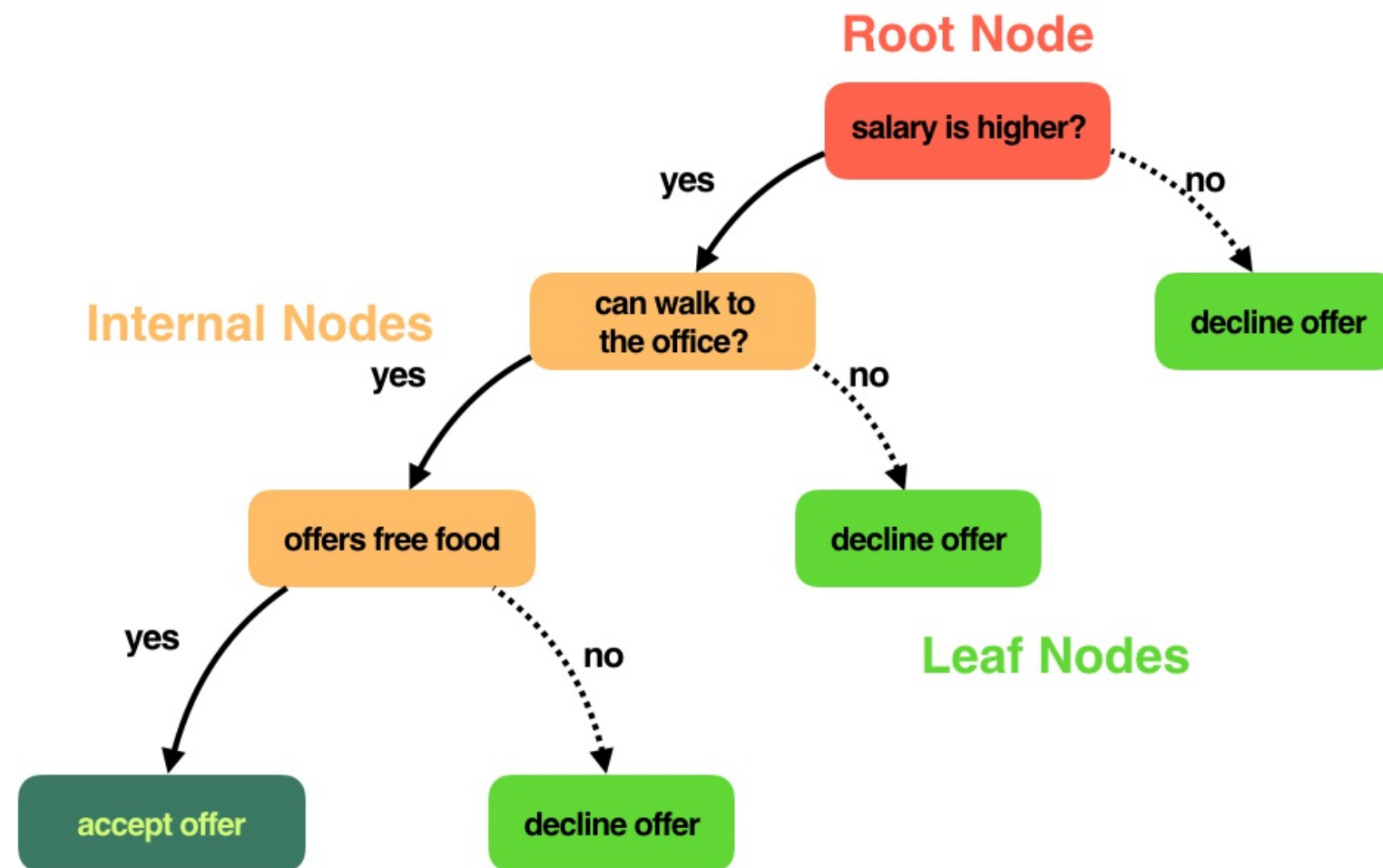


# We will cover:

- Classification & Regression Trees
- Bagged Trees
- Random Forests
- Boosted Trees (GBM)



# Decision tree terminology: nodes





# Training Decision Trees in R

```
> library("rpart")
```

```
> help(package = "rpart")
```

## Recursive Partitioning and Regression Trees



### Documentation for package 'rpart' version 4.1-10

- [DESCRIPTION file.](#)
- [User guides, package vignettes and other documentation.](#)
- [Package NEWS.](#)

### Help Pages

<a href="#">car.test.frame</a>	Automobile Data from 'Consumer Reports' 1990
<a href="#">car90</a>	Automobile Data from 'Consumer Reports' 1990
<a href="#">cu.summary</a>	Automobile Data from 'Consumer Reports' 1990
<a href="#">kyphosis</a>	Data on Children who have had Corrective Spinal Surgery
<a href="#">labels.rpart</a>	Create Split Labels For an Rpart Object
<a href="#">meanvar</a>	Mean-Variance Plot for an Rpart Object
<a href="#">meanvar.rpart</a>	Mean-Variance Plot for an Rpart Object
<a href="#">na.rpart</a>	Handles Missing Values in an Rpart Object
<a href="#">path.rpart</a>	Follow Paths to Selected Nodes of an Rpart Object
<a href="#">plot.rpart</a>	Plot an Rpart Object
<a href="#">plotcp</a>	Plot a Complexity Parameter Table for an Rpart Fit
<a href="#">post</a>	PostScript Presentation Plot of an Rpart Object
<a href="#">post.rpart</a>	PostScript Presentation Plot of an Rpart Object
<a href="#">predict.rpart</a>	Predictions from a Fitted Rpart Object
<a href="#">print.rpart</a>	Print an Rpart Object
<a href="#">printcp</a>	Displays CP table for Fitted Rpart Object
<a href="#">prune</a>	Cost-complexity Pruning of an Rpart Object
<a href="#">prune.roart</a>	Cost-complexity Pruning of an Rpart Object
<a href="#">residuals.roart</a>	Residuals From a Fitted Rpart Object
<a href="#">rpart</a>	Recursive Partitioning and Regression Trees
<a href="#">rpart.control</a>	Control for Rpart Fits
<a href="#">rpart.exp</a>	Initialization function for exponential fitting
<a href="#">rpart.object</a>	Recursive Partitioning and Regression Trees Object
<a href="#">rsq.rpart</a>	Plots the Approximate R-Square for the Different Splits
<a href="#">snip.rpart</a>	Snip Subtrees of an Rpart Object
<a href="#">solder</a>	Soldering of Components on Printed-Circuit Boards
<a href="#">stagec</a>	Stage C Prostate Cancer
<a href="#">summary.rpart</a>	Summarize a Fitted Rpart Object
<a href="#">text.rpart</a>	Place Text on a Dendrogram Plot
<a href="#">xpred.rpart</a>	Return Cross-Validated Predictions



# Training Decision Trees in R

```
> rpart(response ~ ., data = dataset)
```



## MACHINE LEARNING WITH TREE-BASED MODELS IN R

**Let's practice!**





MACHINE LEARNING WITH TREE-BASED MODELS IN R

# Introduction to classification trees

Gabriela de Queiroz  
Instructor



# Advantages

- ✓ Simple to understand, interpret, visualize
- ✓ Can handle both numerical and categorical features (inputs) natively
- ✓ Can handle missing data elegantly
- ✓ Robust to outliers
- ✓ Requires little data preparation
- ✓ Can model non-linearity in the data
- ✓ Can be trained quickly on large datasets



# Disadvantages

- ✖ Large trees can be hard to interpret
- ✖ Trees have high variance, which causes model performance to be poor
- ✖ Trees overfit easily



# Will you wait for a table or go elsewhere?

customer	fri/sat	raining	reservation	wait estimate	will_wait?
1	No	No	Yes	0-10	Yes
2	No	No	No	30-60	No
3	No	No	No	0-10	Yes
4	Yes	No	No	10-30	Yes
5	Yes	No	Yes	> 60	No
6	No	Yes	Yes	0-10	Yes
...	...	...	...	...	...

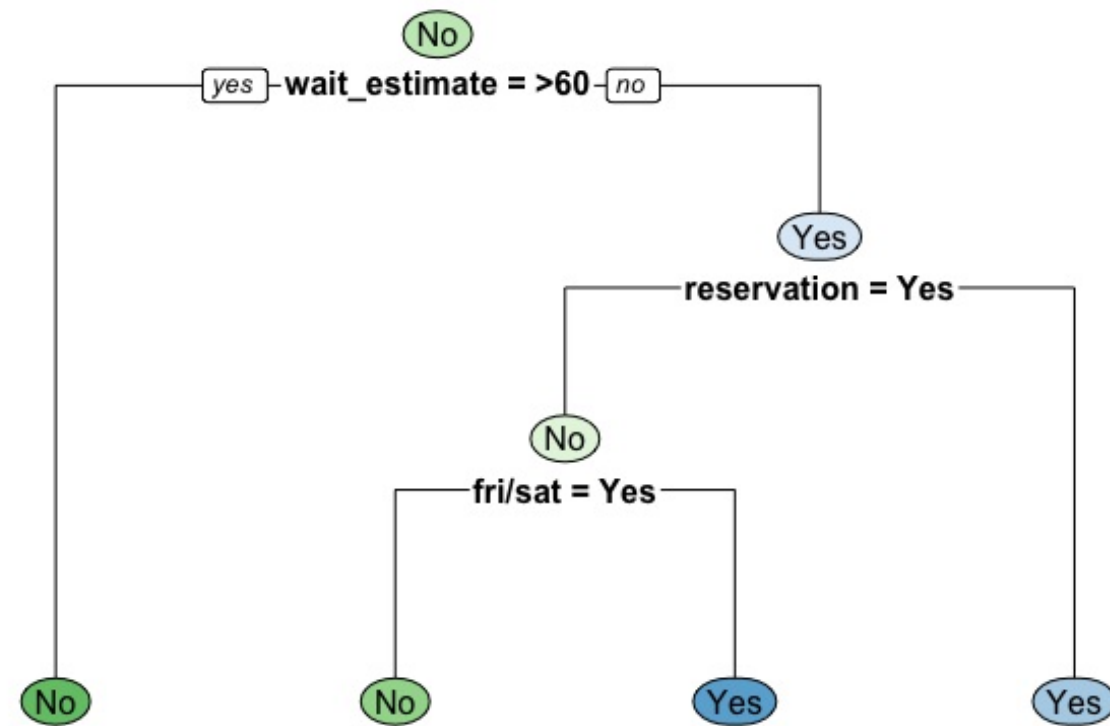


# Restaurant Example

customer	fri/sat	raining	reservation	wait estimate	will_wait?
1	No	No	Yes	0-10	Yes
2	No	No	No	30-60	No
3	No	No	No	0-10	Yes
4	Yes	No	No	10-30	Yes
5	Yes	No	Yes	> 60	No
6	No	Yes	Yes	0-10	Yes
...	...	...	...	...	...



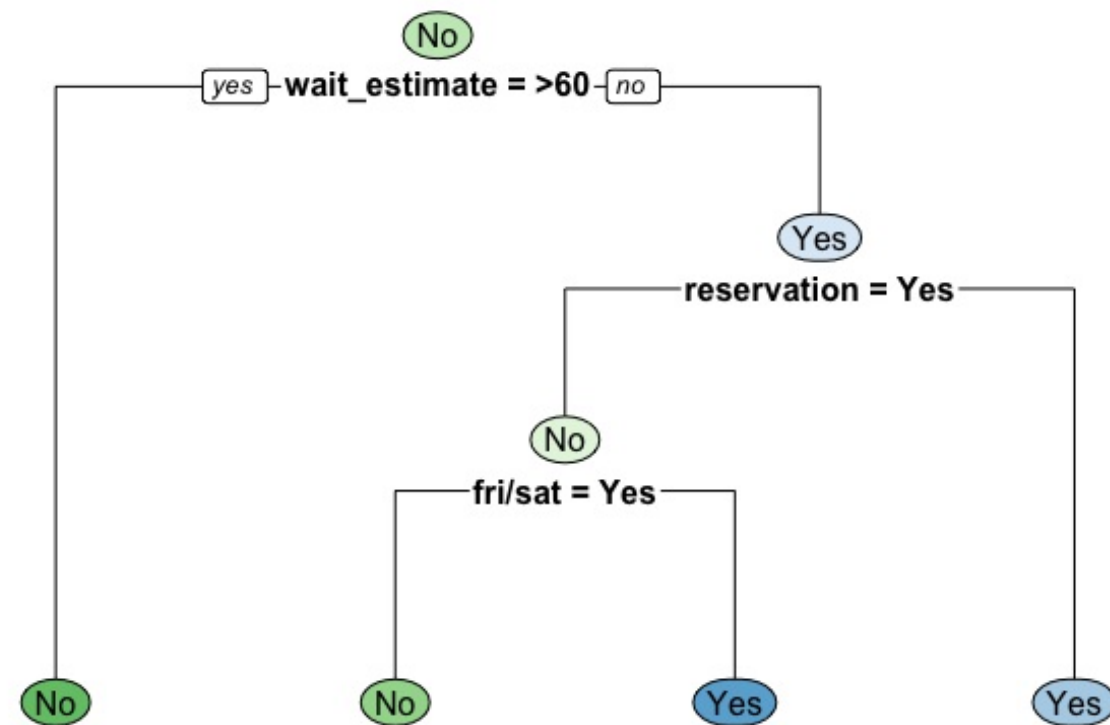
# Decision Tree in R





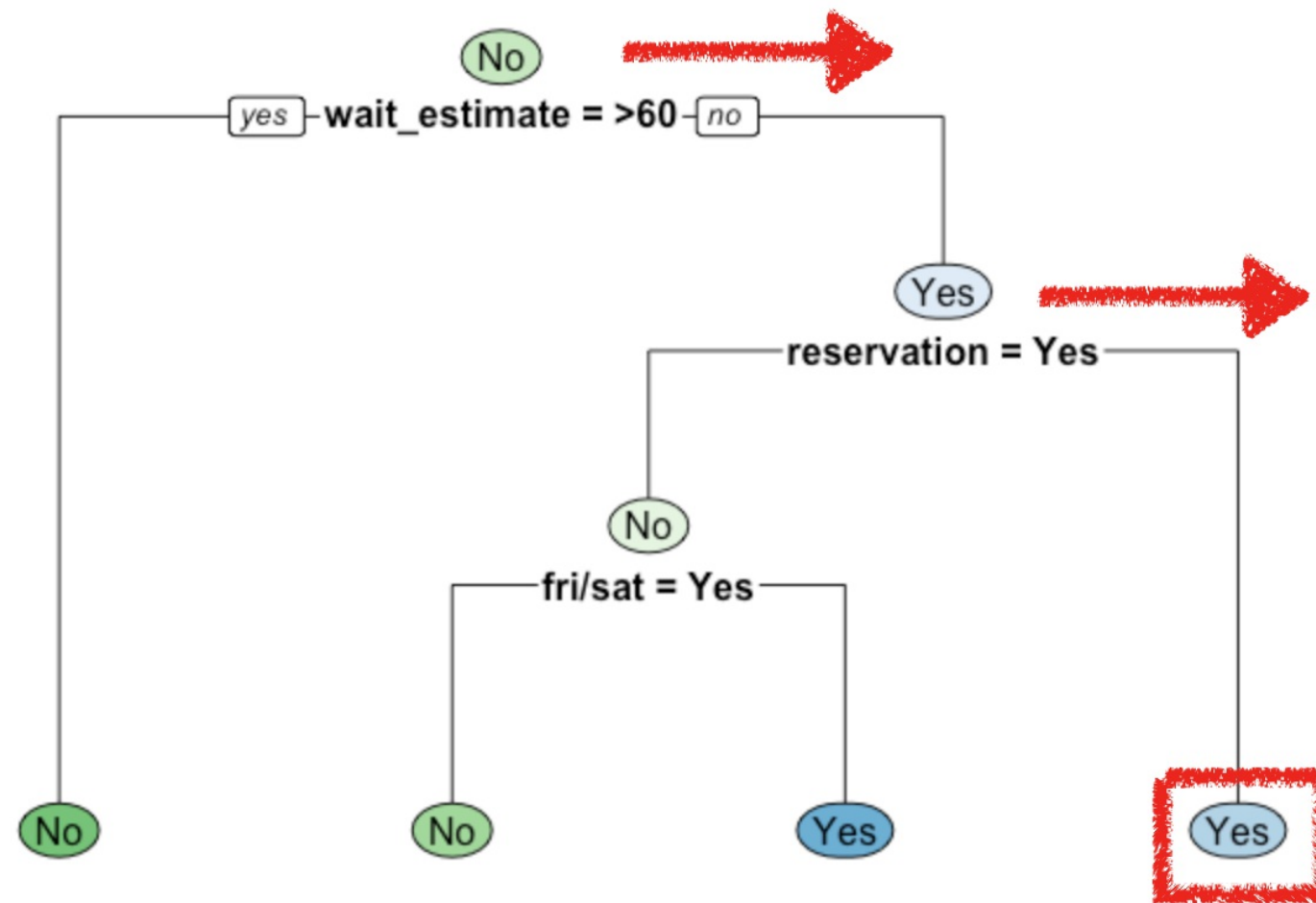
# Prediction example

- The wait estimate is 20 minutes, no reservation was made, and it is Wednesday





# Example







## MACHINE LEARNING WITH TREE-BASED MODELS IN R

**Let's practice!**



MACHINE LEARNING WITH TREE-BASED MODELS IN R

# Overview of the modeling process

Gabriela de Queiroz  
Instructor



# Train/Test Split





# Train/test split in R

```
# total number of rows in the restaurant data frame
n <- nrow(restaurant)

# number of rows for the training set (80% of the dataset)
n_train <- round(0.80 * n)

# create a vector of indices which is an 80% random sample
set.seed(123) # set a random seed for reproducibility
train_indices <- sample(1:n, n_train)

# subset the data frame to training indices only
restaurant_train <- restaurant[train_indices, ]

# exclude the training indices to create the test set
restaurant_test <- restaurant[-train_indices, ]
```



# Train a Classification Tree

```
# train the model to predict the binary response, "will_wait"

restaurant_model <- rpart(formula = will_wait ~.,
                           data = restaurant_train,
                           method = "class")
```

**formula:** response variable ~ predictor variables



## MACHINE LEARNING WITH TREE-BASED MODELS IN R

**Let's practice!**



MACHINE LEARNING WITH TREE-BASED MODELS IN R

# Evaluate Model Performance

Gabriela de Queiroz  
Instructor



# Predicting class labels for test data

```
> predict(model, test_dataset)
```

```
> predict(model, test_dataset, type = ____)
```

```
class_prediction <- predict(object = restaurant_model, # model object  
                             newdata = restaurant_test, # test dataset  
                             type = "class") # return classification labels
```





# Evaluation Metrics for Binary Classification

- Accuracy
- Confusion Matrix
- Log-loss
- AUC



# Accuracy

$$accuracy = \frac{\text{n of correct predictions}}{\text{n of total data points}}$$



# Confusion Matrix

		Actual	
Predicted	YES	YES	NO
	NO	YES	NO



# Confusion Matrix

		Actual	
		YES	NO
Predicted	YES	TRUE POSITIVE (TP)	FALSE POSITIVE (FP)
	NO	FALSE NEGATIVE (FN)	TRUE NEGATIVE (TN)



# Confusion Matrix

```
library(caret)

# calculate the confusion matrix for the test set
confusionMatrix(data = class_prediction,          # predicted classes
                 reference = restaurant_test$will_wait) # actual classes
```



## MACHINE LEARNING WITH TREE-BASED MODELS IN R

**Let's practice!**



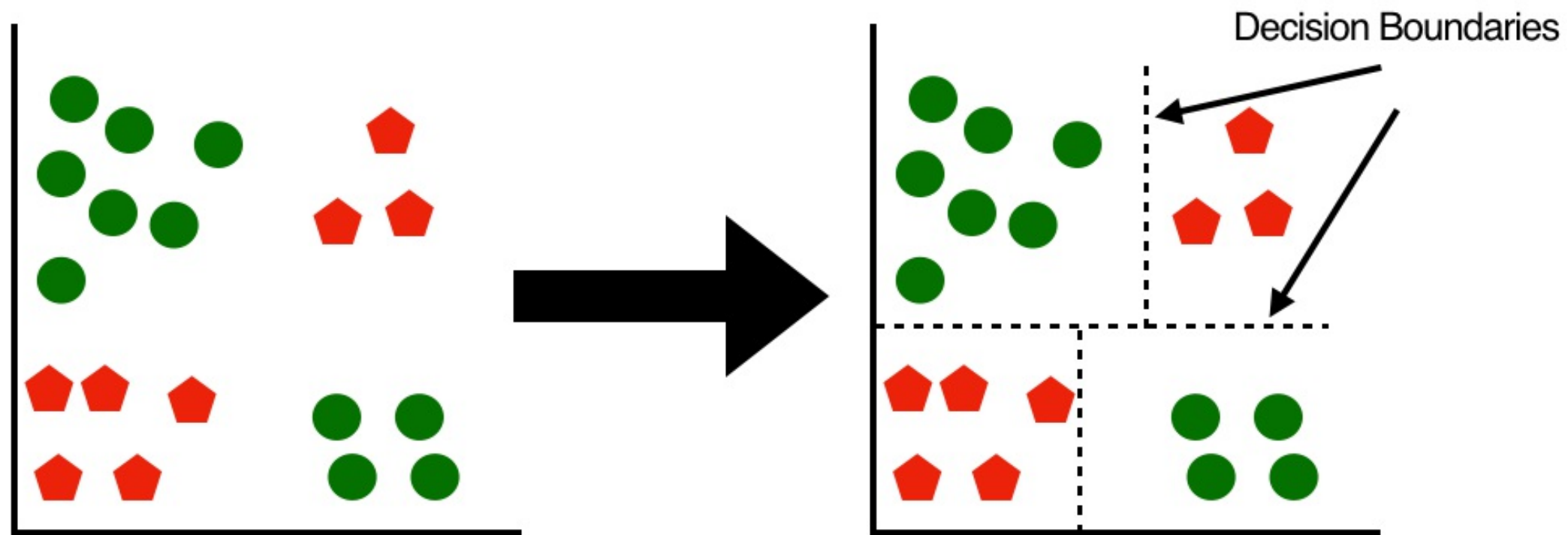
MACHINE LEARNING WITH TREE-BASED MODELS IN R

# Use of splitting criterion in trees

Gabriela de Queiroz  
Instructor



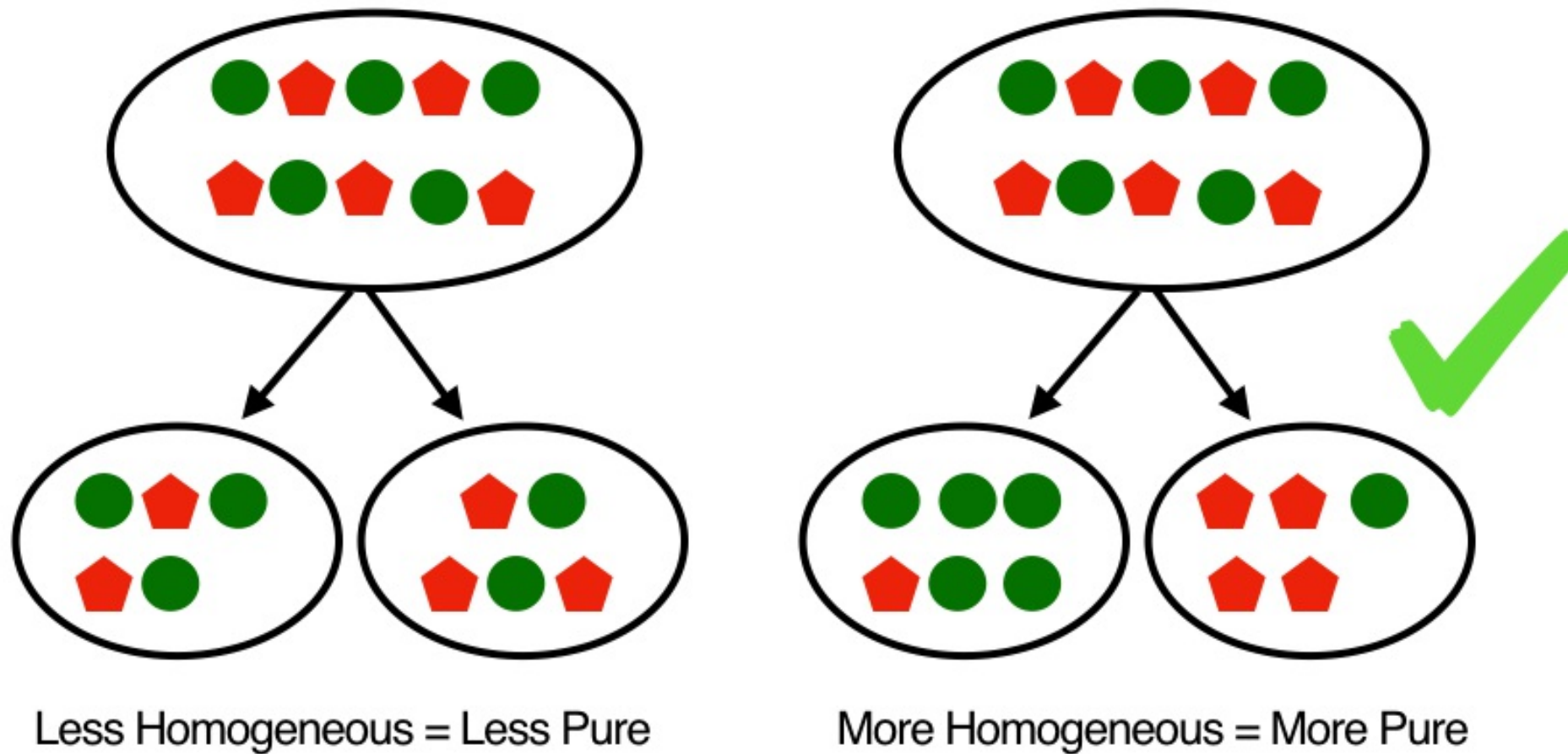
# Split the data into "pure" regions



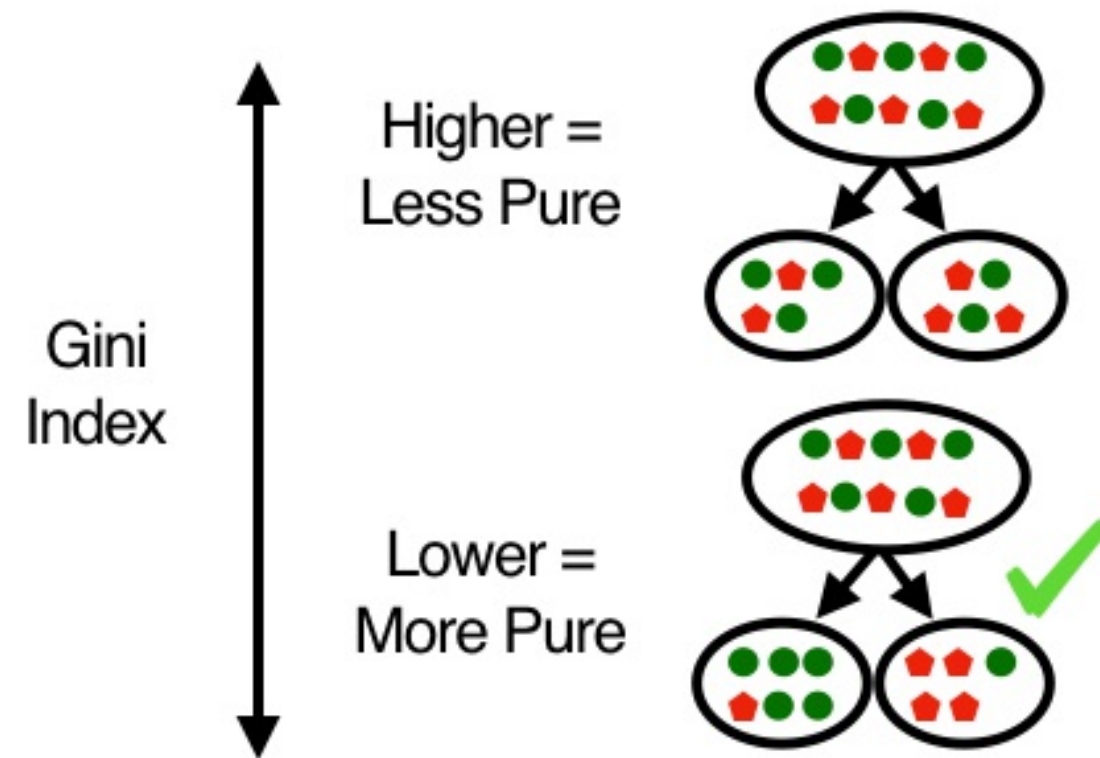




# How to determine the best split?



# Impurity Measure - Gini Index





## MACHINE LEARNING WITH TREE-BASED MODELS IN R

**Let's practice!**