

# Introduction to Machine Learning

## Regression and Classification Trees

Andres Mendez-Vazquez

July 8, 2018

# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



# Decision Trees

Powerful/popular

For classification and prediction.

# Decision Trees

## Powerful/popular

For classification and prediction.

## Represent rules

- Rules can be expressed in English.

▶ *IF Age  $\leq$  43 & Sex == Male AND  
Credit Card Insurance == No THEN*

*Life Insurance Promotion = No*

- Rules can be expressed using SQL for query.

# Decision Trees

## Powerful/popular

For classification and prediction.

## Represent rules

- Rules can be expressed in English.
  - ▶ **IF** *Age*  $\leq$  43 & *Sex* == *Male* **AND**  
*Credit Card Insurance* == *No* **THEN**

*Life Insurance Promotion* = *No*

- Rules can be expressed using SQL for query.

Useful to explore data to gain insight into relationships

Of a large number of candidate input variables to a target (output) variable.

# Decision Trees

## Powerful/popular

For classification and prediction.

## Represent rules

- Rules can be expressed in English.
  - ▶ **IF** *Age*  $\leq$  43 & *Sex* == *Male* **AND**  
*Credit Card Insurance* == *No* **THEN**  
  
*Life Insurance Promotion* = *No*
- Rules can be expressed using SQL for query.

Learn to explore the hidden insights in relationships

Of a large number of candidate input variables to a target (output) variable.

# Decision Trees

## Powerful/popular

For classification and prediction.

## Represent rules

- Rules can be expressed in English.
  - ▶ **IF** *Age*  $\leq$  43 & *Sex* == *Male* **AND**  
*Credit Card Insurance* == *No* **THEN**

*Life Insurance Promotion* = *No*

- Rules can be expressed using SQL for query.

## Useful to explore data to gain insight into relationships

Of a large number of candidate input variables to a target (output) variable.



# What are They?

## Decision Tree

A structure that can be used to divide up a large collection of records into successively smaller sets of records by applying a sequence of simple decision rules.

### A decision tree model

Consists of a set of rules for dividing a large heterogeneous population into smaller, more homogeneous groups with respect to a particular target variable.



# What are They?

## Decision Tree

A structure that can be used to divide up a large collection of records into successively smaller sets of records by applying a sequence of simple decision rules.

## A decision tree model

Consists of a set of rules for dividing a large heterogeneous population into smaller, more homogeneous groups with respect to a particular target variable.



# Decision Tree Types

## Binary trees

- Only two choices in each split. Can be non-uniform (uneven) in depth.

## N-way trees or Ternary Trees

- Three or more choices in at least one of its splits (3-way, 4-way, etc.)



# Decision Tree Types

## Binary trees

- Only two choices in each split. Can be non-uniform (uneven) in depth.

## N-way trees or Ternary trees

- Three or more choices in at least one of its splits (3-way, 4-way, etc.).



# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- **Examples of Trees**

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

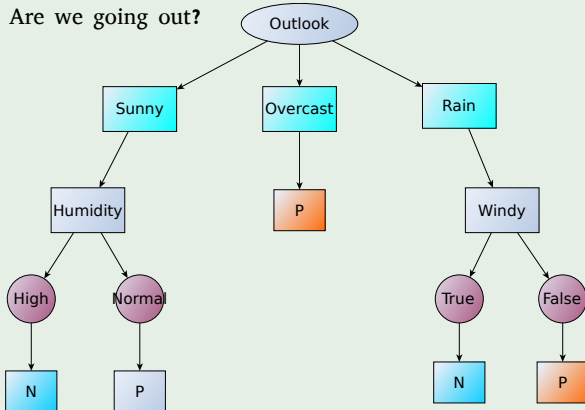
## 5 Conclusions

- First Some Remarks
- Issues



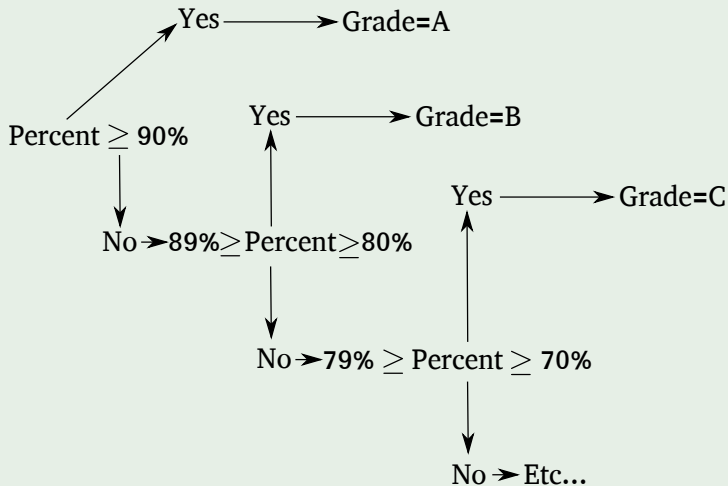
# An Example

We have



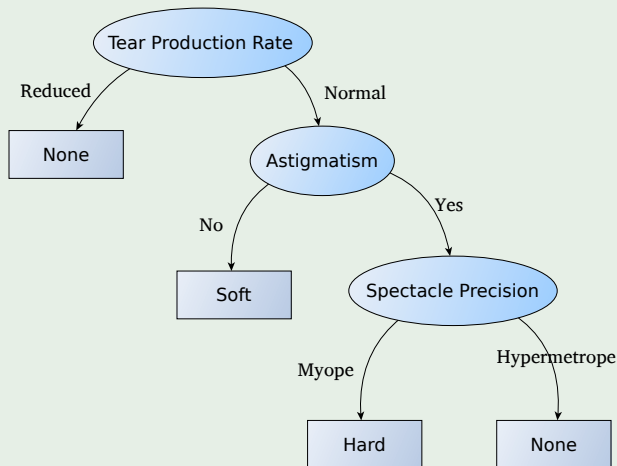
## Another Example - Grades

### Deciding the grades



# Yet Another Example

## Decision About Needing Glasses





# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



# First than anything

## Assume

Consider a Regression Problem with:

- ① Continuous Response  $y$ .
- ② Inputs  $x_1$  and  $x_2$  taking values in  $[0, 1]$ .
- ③ We have only recursive binary decisions/partitions.



# First than anything

## Assume

Consider a Regression Problem with:

- 1 Continuous Response  $y$ .
- 2 Inputs  $x_1$  and  $x_2$  taking values in  $[0, 1]$ .
- 3 We have only recursive binary decisions/partitions.

Example of a partition



# First than anything

## Assume

Consider a Regression Problem with:

- 1 Continuous Response  $y$ .
- 2 Inputs  $x_1$  and  $x_2$  taking values in  $[0, 1]$ .

⊙ We have only recursive binary decisions/partitions.

Example of a partition



# First than anything

## Assume

Consider a Regression Problem with:

- 1 Continuous Response  $y$ .
- 2 Inputs  $x_1$  and  $x_2$  taking values in  $[0, 1]$ .
- 3 We have only recursive binary decisions/partitions.

Example of a partition



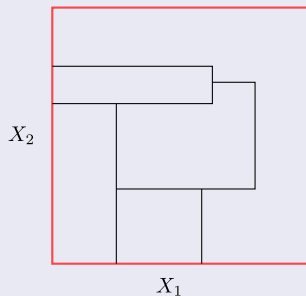
# First than anything

## Assume

Consider a Regression Problem with:

- 1 Continuous Response  $y$ .
- 2 Inputs  $x_1$  and  $x_2$  taking values in  $[0, 1]$ .
- 3 We have only recursive binary decisions/partitions.

## Example of a partition



# Everything is fine, but!!!

## Although

- In each partition element we can model  $Y$  with a different constant.



# Everything is fine, but!!!

## Although

- In each partition element we can model  $Y$  with a different constant.

## There is a problem

- Each partitioning line has a simple description like  $x_1 = c!!!$
- The Resulting Regions are difficult to describe!!!





# Everything is fine, but!!!

## Although

- In each partition element we can model  $Y$  with a different constant.

## There is a problem

- Each partitioning line has a simple description like  $x_1 = c!!!$
- The Resulting Regions are difficult to describe!!!



# Everything is fine, but!!!

## Although

- In each partition element we can model  $Y$  with a different constant.

## There is a problem

- Each partitioning line has a simple description like  $x_1 = c!!!$
- The Resulting Regions are difficult to describe!!!



# Solving the Issue

## We do the following

- Chose a variable and split the space using  $x_i = c$

Keep doing that using one of the variables until a rules stops the process

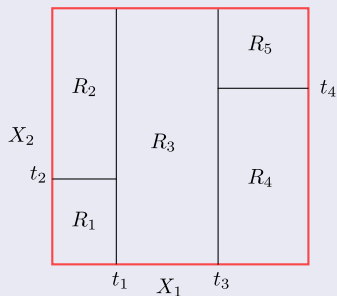


# Solving the Issue

We do the following

- Chose a variable and split the space using  $x_i = c$

Keep doing that using one of the variables until a rules stops the process



# The corresponding Regression Tree

We have

$$\hat{y} = f(\mathbf{x}) = \sum_{m=1}^5 c_m I\{(x_1, x_2) \in R_m\}$$

This regression can be interpreted as

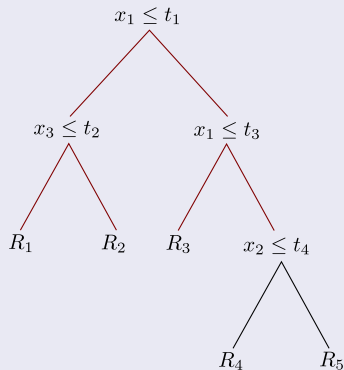


## The corresponding Regression Tree

We have

$$\hat{y} = f(\mathbf{x}) = \sum_{m=1}^5 c_m I\{(x_1, x_2) \in R_m\}$$

This regression can be interpreted as



# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- **Structure of Decision Trees**
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



# Structure

## Structure

- Nodes
  - ▶ Appear as rectangles or circles
    - ▶ Represent test or decision
  - Lines or branches - represent outcome of a test
  - Circles - terminal (leaf) nodes.





# Structure

## Structure

- Nodes
  - ▶ Appear as rectangles or circles
  - ▶ Represent test or decision
- Lines or branches - represent outcome of a test
- Circles - terminal (leaf) nodes.

## Nodes

- Top or starting node is root node
- Internal nodes are used for decisions
- Terminal Nodes or Leaves are the final results



# Structure

## Structure

- Nodes
  - ▶ Appear as rectangles or circles
  - ▶ Represent test or decision
- Lines or branches - represent outcome of a test
- Circles - terminal (leaf) nodes.

## Nodes

- Top or starting node is root node
- Internal nodes are used for decisions
- Terminal Nodes or Leaves are the final results



# Structure

## Structure

- Nodes
  - ▶ Appear as rectangles or circles
  - ▶ Represent test or decision
- Lines or branches - represent outcome of a test
- Circles - terminal (leaf) nodes.

## Nodes

- Top or starting node is root node
- Internal nodes are used for decisions
- Terminal Nodes or Leaves are the final results



# Structure

## Structure

- Nodes
  - ▶ Appear as rectangles or circles
  - ▶ Represent test or decision
- Lines or branches - represent outcome of a test
- Circles - terminal (leaf) nodes.

## Nodes

- Top or starting node is root node

• Internal nodes are used for decisions

• Terminal Nodes or Leaves are the final results



# Structure

## Structure

- Nodes
  - ▶ Appear as rectangles or circles
  - ▶ Represent test or decision
- Lines or branches - represent outcome of a test
- Circles - terminal (leaf) nodes.

## Nodes

- Top or starting node is root node
- Internal nodes are used for decisions

• Terminal Nodes or Leaves are the final results



# Structure

## Structure

- Nodes
  - ▶ Appear as rectangles or circles
  - ▶ Represent test or decision
- Lines or branches - represent outcome of a test
- Circles - terminal (leaf) nodes.

## Nodes

- Top or starting node is root node
- Internal nodes are used for decisions
- Terminal Nodes or Leaves are the final results



# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- **Types of Decision Trees**

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



# Types of Decision Trees

## Regression Trees

The predicted outcome can be considered a number.

## Classification Trees

- The predicted outcome is the class to which the data belongs.





# Types of Decision Trees

## Regression Trees

The predicted outcome can be considered a number.

## Classification Trees

- The predicted outcome is the class to which the data belongs.



# Classification and Regression Trees (CART)

## CART

- The term CART is an umbrella term used to refer to both of the above procedures.

### introduced by

- It was introduced by Breiman et. al in the book
  - ▶ “Classification and Regression Trees”

### Similarities

- Regression and Classification trees have some similarities – nevertheless they differ in the way the splitting at each node is done.



# Classification and Regression Trees (CART)

## CART

- The term CART is an umbrella term used to refer to both of the above procedures.

## Introduced by

- It was introduced by Breiman et. al in the book
  - ▶ **“Classification and Regression Trees”**

## Similarities

- Regression and Classification trees have some similarities – nevertheless they differ in the way the splitting at each node is done.



# Classification and Regression Trees (CART)

## CART

- The term CART is an umbrella term used to refer to both of the above procedures.

## Introduced by

- It was introduced by Breiman et. al in the book
  - ▶ **“Classification and Regression Trees”**

## Similarities

- Regression and Classification trees have some similarities – nevertheless they differ in the way the splitting at each node is done.



# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- **Growing Regression Trees**
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



# Setup

Data Consists on inputs of dimensionality  $d$

$$\left\{ (x_i, y_i)_{i=1}^N \right\}$$

Where  $x_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$ .

Here we want an algorithm

- To do the splitting automatically

This assume a final  $M$  partition  $R_1, R_2, \dots, R_M$

- We model the response as a constant  $c_m$  in each region

$$f(\mathbf{x}) = \sum_{m=1}^M c_m I(\mathbf{x} \in R_m)$$

# Setup

Data Consists on inputs of dimensionality  $d$

$$\left\{ (x_i, y_i)_{i=1}^N \right\}$$

Where  $x_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$ .

Here, we want an algorithm

- To do the splitting automatically

What is the optimal way to partition  $R \subset \mathbb{R}^d$  into  $M$  regions?

- We model the response as a constant  $c_m$  in each region

$$f(\mathbf{x}) = \sum_{m=1}^M c_m I(\mathbf{x} \in R_m)$$

# Setup

Data Consists on inputs of dimensionality  $d$

$$\left\{ (x_i, y_i)_{i=1}^N \right\}$$

Where  $x_i = (x_{i1}, x_{i2}, \dots, x_{id})^T$ .

Here, we want an algorithm

- To do the splitting automatically

Thus, assume a initial  $M$  partition  $R_1, R_2, \dots, R_M$

- We model the response as a constant  $c_m$  in each region

$$f(\mathbf{x}) = \sum_{m=1}^M c_m I(\mathbf{x} \in R_m)$$



# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- **Using the Sum of Squared Error**
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



We have then

We adopt as our criterion minimization

$$L(c_1, c_2, \dots, c_M) = \sum_{i=1}^N \sum_{m=1}^M (y_i - f(\mathbf{x}_i))^2$$

Then using a classic derivative with respect to

$$\frac{\partial L(c_1, c_2, \dots, c_M)}{\partial c_m} = -2 \sum_{i=1}^N \left( y_i - \sum_{m=1}^M c_m I(\mathbf{x}_i \in R_m) \right) I(\mathbf{x}_i \in R_m)$$

Then

$$\sum_{y_i | \mathbf{x}_i \in R_m} y_i - \sum_{i=1}^N I(\mathbf{x}_i \in R_m) \sum_{m=1}^M c_m I(\mathbf{x}_i \in R_m) = 0$$

We have then

We adopt as our criterion minimization

$$L(c_1, c_2, \dots, c_M) = \sum_{i=1}^N \sum_{m=1}^M (y_i - f(\mathbf{x}_i))^2$$

Then using a classic derivative with respect to  $c_m$

$$\frac{\partial L(c_1, c_2, \dots, c_M)}{\partial c_m} = -2 \sum_{i=1}^N \left( y_i - \sum_{m=1}^M c_m I(\mathbf{x}_i \in R_m) \right) I(\mathbf{x}_i \in R_m)$$

Then

$$\sum_{y_i | \mathbf{x}_i \in R_m} y_i - \sum_{i=1}^N I(\mathbf{x}_i \in R_m) \sum_{m=1}^M c_m I(\mathbf{x}_i \in R_m) = 0$$

We have then

We adopt as our criterion minimization

$$L(c_1, c_2, \dots, c_M) = \sum_{i=1}^N \sum_{m=1}^M (y_i - f(\mathbf{x}_i))^2$$

Then using a classic derivative with respect to  $c_m$

$$\frac{\partial L(c_1, c_2, \dots, c_M)}{\partial c_m} = -2 \sum_{i=1}^N \left( y_i - \sum_{m=1}^M c_m I(\mathbf{x}_i \in R_m) \right) I(\mathbf{x}_i \in R_m)$$

Then

$$\sum_{y_i | \mathbf{x}_i \in R_m} y_i - \sum_{i=1}^N I(\mathbf{x}_i \in R_m) \sum_{m=1}^M c_m I(\mathbf{x}_i \in R_m) = 0$$

# The simplest function for $c_m$

## Something Notable

$$\sum_{\mathbf{x}_i \in R_m} c_m = \sum_{y_i | \mathbf{x}_i \in R_m} y_i$$

Then

$$c_m = \frac{1}{N_m} \sum_{y_i | \mathbf{x}_i \in R_m} y_i$$

Problem:

- Finding the best binary partition in terms of minimum sum of squares is generally  $O(2^N)$  a NP Problem!!!



# The simplest function for $c_m$

## Something Notable

$$\sum_{\mathbf{x}_i \in R_m} c_m = \sum_{y_i | \mathbf{x}_i \in R_m} y_i$$

## Then

$$c_m = \frac{1}{N_m} \sum_{y_i | \mathbf{x}_i \in R_m} y_i$$

## Problem:

- Finding the best binary partition in terms of minimum sum of squares is generally  $O(2^N)$  a NP Problem!!!



# The simplest function for $c_m$

## Something Notable

$$\sum_{\mathbf{x}_i \in R_m} c_m = \sum_{y_i | \mathbf{x}_i \in R_m} y_i$$

## Then

$$c_m = \frac{1}{N_m} \sum_{y_i | \mathbf{x}_i \in R_m} y_i$$

## Problem

- Finding the best binary partition in terms of minimum sum of squares is generally  $O(2^N)$  a NP Problem!!!



# What to do?

Consider a splitting variable  $j$  and split point  $s$

- Define the pair of half-planes

$$R_1(j, s) = \{\mathbf{x} | x_j \leq s\} \text{ and } R_2(j, s) = \{\mathbf{x} | x_j > s\}$$

Using an Optimization Problem

$$\min_{j,s} \left\{ \min_{c_1} \sum_{\mathbf{x}_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{\mathbf{x}_i \in R_2(j,s)} (y_i - c_2)^2 \right\}$$

The nice part of this

- For any choice  $j$  and  $s$ , the inner minimization is solved by

$$\hat{c}_1 = \frac{1}{N_1} \sum_{y_i | \mathbf{x}_i \in R_1(j,s)} y_i \text{ and } \hat{c}_2 = \frac{1}{N_2} \sum_{y_i | \mathbf{x}_i \in R_2(j,s)} y_i$$



# What to do?

Consider a splitting variable  $j$  and split point  $s$

- Define the pair of half-planes

$$R_1(j, s) = \{\mathbf{x} | x_j \leq s\} \text{ and } R_2(j, s) = \{\mathbf{x} | x_j > s\}$$

Using an Optimization Problem

$$\min_{j,s} \left\{ \min_{c_1} \sum_{\mathbf{x}_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{\mathbf{x}_i \in R_2(j,s)} (y_i - c_2)^2 \right\}$$

The nice part of this

- For any choice  $j$  and  $s$ , the inner minimization is solved by

$$\hat{c}_1 = \frac{1}{N_1} \sum_{\mathbf{x}_i \in R_1(j,s)} y_i \text{ and } \hat{c}_2 = \frac{1}{N_2} \sum_{\mathbf{x}_i \in R_2(j,s)} y_i$$

# What to do?

Consider a splitting variable  $j$  and split point  $s$

- Define the pair of half-planes

$$R_1(j, s) = \{\mathbf{x} | x_j \leq s\} \text{ and } R_2(j, s) = \{\mathbf{x} | x_j > s\}$$

Using an Optimization Problem

$$\min_{j,s} \left\{ \min_{c_1} \sum_{\mathbf{x}_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{\mathbf{x}_i \in R_2(j,s)} (y_i - c_2)^2 \right\}$$

The nice part of this

- For any choice  $j$  and  $s$ , the inner minimization is solved by

$$\hat{c}_1 = \frac{1}{N_1} \sum_{y_i | \mathbf{x}_i \in R_1(j,s)} y_i \text{ and } \hat{c}_2 = \frac{1}{N_2} \sum_{y_i | \mathbf{x}_i \in R_2(j,s)} y_i$$

# Therefore

For each splitting variable  $j$

- Finding  $s$  is done quickly!!!

We can repeat this process

- Problem, we can finish with an over-fitting tree/a very large tree.

How do we solve?

- Tree size is an hyper-parameter governing the model's complexity.



Cinvestav

# Therefore

For each splitting variable  $j$

- Finding  $s$  is done quickly!!!

We can repeat this process

- Problem, we can finish with an over-fitting tree/a very large tree.

How do we solve?

- Tree size is an hyper-parameter governing the model's complexity.



# Therefore

For each splitting variable  $j$

- Finding  $s$  is done quickly!!!

We can repeat this process

- Problem, we can finish with an over-fitting tree/a very large tree.

How do we solve?

- Tree size is an hyper-parameter governing the model's complexity.



# Therefore

We have that

- Tree size is a tuning parameter governing the model's complexity

A preferred strategy:

- Grow the tree until some minimum size node is done.

Then

- This large tree is pruned using cost-complexity pruning.



# Therefore

## We have that

- Tree size is a tuning parameter governing the model's complexity

## A preferred strategy

- Grow the tree until some minimum size node is done.

## Then

- This large tree is pruned using cost-complexity pruning.



# Therefore

## We have that

- Tree size is a tuning parameter governing the model's complexity

## A preferred strategy

- Grow the tree until some minimum size node is done.

## Then

- This large tree is pruned using cost-complexity pruning.





# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- **Pruning**

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



# We need to define something

## Definition

- We define a subtree  $T \subseteq T_0$  to be any tree that can be obtained by pruning  $T_0$ :
  - ▶ By collapsing any number of its internal (non-terminal) nodes.

Given that each  $T_0$  is indexed by  $m$

- Let  $|T|$  denote the number of terminal nodes in  $T$ :

$$N_m = |R_m|, \quad \bar{c}_m = \frac{1}{N_m} \sum_{y_i | x_i \in R_m} y_i \quad \text{and} \quad Q_m(T) = \frac{1}{N_m} (\bar{c}_m - y_i)^2$$



# We need to define something

## Definition

- We define a subtree  $T \subseteq T_0$  to be any tree that can be obtained by pruning  $T_0$ :
  - ▶ By collapsing any number of its internal (non-terminal) nodes.

## Given that each $R_m$ is indexed by $m$

- Let  $|T|$  denote the number of terminal nodes in  $T$ :

$$N_m = |R_m|, \hat{c}_m = \frac{1}{N_m} \sum_{y_i | x_i \in R_m} y_i \text{ and } Q_m(T) = \frac{1}{N_m} (\hat{c}_m - y_i)^2$$



Thus

Define the cost complexity criterion with  $\alpha \geq 0$

$$C_\alpha(T) = \sum_{m=1}^{|T|} N_m Q_m(T) + \alpha |T|$$

Finally

- The idea is to find, for each  $\alpha$ , the subtree  $T_\alpha \subseteq T_0$  to minimize  $C_\alpha(T)$

Properties of  $\alpha$

- Large values of  $\alpha$  result in smaller  $T_\alpha$
- Small values of  $\alpha$  result in larger  $T_\alpha$



Thus

Define the cost complexity criterion with  $\alpha \geq 0$

$$C_\alpha(T) = \sum_{m=1}^{|T|} N_m Q_m(T) + \alpha |T|$$

Finally

- The idea is to find, for each  $\alpha$ , the subtree  $T_\alpha \subseteq T_0$  to minimize  $C_\alpha(T)$

Properties of  $T_\alpha$

- Large values of  $\alpha$  result in smaller  $T_\alpha$
- Small values of  $\alpha$  result in larger  $T_\alpha$



Thus

Define the cost complexity criterion with  $\alpha \geq 0$

$$C_\alpha(T) = \sum_{m=1}^{|T|} N_m Q_m(T) + \alpha |T|$$

Finally

- The idea is to find, for each  $\alpha$ , the subtree  $T_\alpha \subseteq T_0$  to minimize  $C_\alpha(T)$

Properties of  $\alpha$

- Large values of  $\alpha$  result in smaller  $T_\alpha$
- Small values of  $\alpha$  result in larger  $T_\alpha$



## Furthermore

For each  $\alpha$  one can show the existence of unique smallest subtree  $T_\alpha$

- How do we find  $T_\alpha$ ?

Simplest link pruning

- We successively collapse the internal node that produces the smallest per-node increase in

$$\sum_{m=1}^{|T|} N_m Q_m (T)$$

Until you get a single-node (root) and a sequence

$$T \supseteq T_1 \supseteq T_2 \supseteq \dots \supseteq T_N$$

## Furthermore

For each  $\alpha$  one can show the existence of unique smallest subtree  $T_\alpha$

- How do we find  $T_\alpha$ ?

## Using weakest link pruning

- We successively collapse the internal node that produces the smallest per-node increase in

$$\sum_{m=1}^{|T|} N_m Q_m (T)$$

Until you get a single-node (root) and a sequence

$$T \supseteq T_1 \supseteq T_2 \supseteq \dots \supseteq T_N$$



## Furthermore

For each  $\alpha$  one can show the existence of unique smallest subtree  $T_\alpha$

- How do we find  $T_\alpha$ ?

## Using weakest link pruning

- We successively collapse the internal node that produces the smallest per-node increase in

$$\sum_{m=1}^{|T|} N_m Q_m (T)$$

Until you get a single-node (root) and a sequence

$$T \supseteq T_1 \supseteq T_2 \supseteq \cdots \supseteq T_N$$

# Then

We get that

- $T_\alpha$  is one of the threes in the in the sequence.

Estimation of  $\hat{\alpha}$  is achieved by cross-validation

- We choose the value  $\hat{\alpha}$  to minimize the cross-validated sum of squares.
  - ▶ This is the final  $T_{\hat{\alpha}}$

For Details

- "Pattern Recognition and Neural Networks" by Brian D. Ripley



Cinvestav

# Then

We get that

- $T_\alpha$  is one of the threes in the in the sequence.

Estimation of  $\alpha$  is achieved by cross-validation

- We choose the value  $\hat{\alpha}$  to minimize the cross-validated sum of squares.
  - ▶ This is the final  $T_{\hat{\alpha}}$

For Details

- "Pattern Recognition and Neural Networks" by Brian D. Ripley



Cinvestav

# Then

## We get that

- $T_\alpha$  is one of the threes in the in the sequence.

## Estimation of $\alpha$ is achieved by cross-validation

- We choose the value  $\hat{\alpha}$  to minimize the cross-validated sum of squares.
  - ▶ This is the final  $T_{\hat{\alpha}}$

## For Details

- “Pattern Recognition and Neural Networks” by Brian D. Ripley



Cinvestav

# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- **Definition**
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



# Important

## Most of the work

It focuses on deciding which property test or query should be performed at the node!!!

If the data test is numerical in nature

There is a way to visualize the decision boundaries produced by the decision trees.



# Important

## Most of the work

It focuses on deciding which property test or query should be performed at the node!!!

## If the data test is numerical in nature

There is a way to visualize the decision boundaries produced by the decision trees.



# Definition OBCT

## Definition

They are binary decision trees where the basic question is  $x_i \leq a_i$ ?

## Example



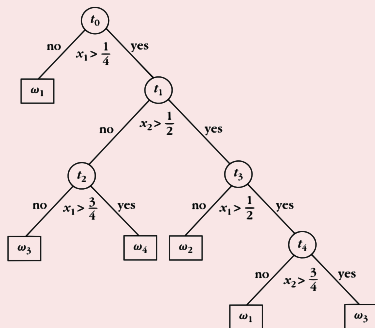
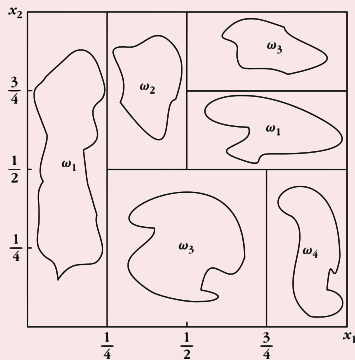


# Definition OBCT

## Definition

They are binary decision trees where the basic question is  $x_i \leq a_i$ ?

## Example



# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- **Training**
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



# Training of a OBCT

## We need first

- At each node, the set of candidate questions to be asked has to be decided.
- Each question corresponds to a specific binary split into two descendant nodes.
- Each node,  $t$ , is associated with a specific subset  $X_t$  of the training set  $X$ .



# Training of a OBCT

## We need first

- At each node, the set of candidate questions to be asked has to be decided.
- Each question corresponds to a specific binary split into two descendant nodes.
- Each node,  $t$ , is associated with a specific subset  $X_t$  of the training set  $X$ .



# Training of a OBCT

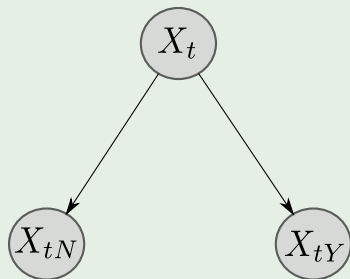
## We need first

- At each node, the set of candidate questions to be asked has to be decided.
- Each question corresponds to a specific binary split into two descendant nodes.
- Each node,  $t$ , is associated with a specific subset  $X_t$  of the training set  $X$ .



## Splitting the Node $X_t$

Basically, we want to split the node into two groups with questions  $t_Y == \text{"YES"}$  and  $t_N = \text{"NO"}$

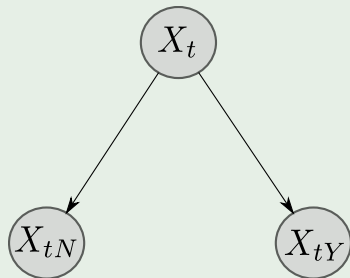


With Properties

- $X_{tY} \cap X_{tN} = \emptyset$ .
- $X_{tY} \cup X_{tN} = X_t$

## Splitting the Node $X_t$

Basically, we want to split the node into two groups with questions  $t_Y == \text{"YES"}$  and  $t_N = \text{"NO"}$



### With Properties

- $X_{tY} \cap X_{tN} = \emptyset$ .
- $X_{tY} \cup X_{tN} = X_t$

# Important

Given the question for each feature  $k$  "Is  $x_k \leq \alpha$ "

For each feature, every possible value of the threshold  $\alpha$  defines a specific split of the subset  $X_t$ .

This in theory

An infinite set of questions has to be asked if  $\alpha$  is an interval  $Y_\alpha \subseteq \mathbb{R}$ .

In practice

only a finite set of questions can be considered.





# Important

Given the question for each feature  $k$  “Is  $x_k \leq \alpha$ ”

For each feature, every possible value of the threshold  $\alpha$  defines a specific split of the subset  $X_t$ .

Thus in theory

An infinite set of questions has to be asked if  $\alpha$  is an interval  $Y_\alpha \subseteq \mathbb{R}$ .

In practice

only a finite set of questions can be considered.



# Important

Given the question for each feature  $k$  “Is  $x_k \leq \alpha$ ”

For each feature, every possible value of the threshold  $\alpha$  defines a specific split of the subset  $X_t$ .

Thus in theory

An infinite set of questions has to be asked if  $\alpha$  is an interval  $Y_\alpha \subseteq \mathbb{R}$ .

In practice

only a finite set of questions can be considered.



## For example

Since the number,  $N$ , of training points in  $X$  is finite

Any of the features  $x_k$  with  $k = 1, \dots, l$  can take at most  $N_t \leq N$  different values

Where

$N_t = |X_t|$  with  $X_t \subset X$

Then

For feature  $x_k$ , one can use  $\alpha_{kn}$  with  $n = 1, 2, \dots, N_{tk}$  and  $N_{tk} \leq N_t$  where  $\alpha_{kn}$  are taken halfway between consecutive distinct values of  $x_k$  in the training subset  $X_t$ .



## For example

Since the number,  $N$ , of training points in  $X$  is finite

Any of the features  $x_k$  with  $k = 1, \dots, l$  can take at most  $N_t \leq N$  different values

Where

$N_t = |X_t|$  with  $X_t \subset X$

When

For feature  $x_k$ , one can use  $\alpha_{kn}$  with  $n = 1, 2, \dots, N_{tk}$  and  $N_{tk} \leq N_t$  where  $\alpha_{kn}$  are taken halfway between consecutive distinct values of  $x_k$  in the training subset  $X_t$ .



## For example

Since the number,  $N$ , of training points in  $X$  is finite

Any of the features  $x_k$  with  $k = 1, \dots, l$  can take at most  $N_t \leq N$  different values

Where

$N_t = |X_t|$  with  $X_t \subset X$

Then

For feature  $x_k$ , one can use  $\alpha_{kn}$  with  $n = 1, 2, \dots, N_{tk}$  and  $N_{tk} \leq N_t$  where  $\alpha_{kn}$  are taken halfway between consecutive distinct values of  $x_k$  in the training subset  $X_t$ .



Then

We repeat this with all features

In such a case, the total number of candidate questions is

$$\sum_{k=1}^l N_{tk} \quad (1)$$

## Then

We repeat this with all features

In such a case, the total number of candidate questions is

$$\sum_{k=1}^l N_{tk} \quad (1)$$

## However

Only one of them has to be chosen to provide the binary split at the current node,  $t$ , of the tree.

## Then

We repeat this with all features

In such a case, the total number of candidate questions is

$$\sum_{k=1}^l N_{tk} \quad (1)$$

## However

Only one of them has to be chosen to provide the binary split at the current node,  $t$ , of the tree.

## Thus

- This is selected to be the one that leads to the best split of the associated subset  $X_t$ .

• The best split is decided according to a splitting criterion.



## Then

We repeat this with all features

In such a case, the total number of candidate questions is

$$\sum_{k=1}^l N_{tk} \quad (1)$$

## However

Only one of them has to be chosen to provide the binary split at the current node,  $t$ , of the tree.

## Thus

- This is selected to be the one that leads to the best split of the associated subset  $X_t$ .
- The best split is decided according to a **splitting criterion**.

# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- **The Sought Criterion**
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



# Criterion's to be Found

## Splitting criterion

- A splitting criterion must be adopted according to which the best split from the set of candidate ones is chosen.

## Stop-splitting rule

A stop-splitting rule is required that controls the growth of the tree, and a node is declared as a terminal one (leaf).

## Rule

A rule is required that assigns each leaf to a specific class.



# Criterion's to be Found

## Splitting criterion

- A splitting criterion must be adopted according to which the best split from the set of candidate ones is chosen.

## Stop-splitting rule

A stop-splitting rule is required that controls the growth of the tree, and a node is declared as a terminal one (leaf).

## Rule

A rule is required that assigns each leaf to a specific class.



# Criterion's to be Found

## Splitting criterion

- A splitting criterion must be adopted according to which the best split from the set of candidate ones is chosen.

## Stop-splitting rule

A stop-splitting rule is required that controls the growth of the tree, and a node is declared as a terminal one (leaf).

## Rule

A rule is required that assigns each leaf to a specific class.



# Looking for Homogeneity!!!

## In order for the tree growing methodology

From the root node down to the leaves every split must generate a subsets that are more homogeneous compared to the ancestor's subset  $X_t$ .

### Meaning

The training feature vectors in each one of the new subsets show, whereas data in  $X_t$  are more equally distributed among the classes.

### Some example

Consider the task of classifying four classes  $\{\omega_1, \omega_2, \omega_3, \omega_4\}$  and assume that the vectors in subset  $X_t$  are distributed among the classes with equal probability.



# Looking for Homogeneity!!!

## In order for the tree growing methodology

From the root node down to the leaves every split must generate a subsets that are more homogeneous compared to the ancestor's subset  $X_t$ .

## Meaning

The training feature vectors in each one of the new subsets show, whereas data in  $X_t$  are more equally distributed among the classes.

Consider the task of classifying four classes  $\{\omega_1, \omega_2, \omega_3, \omega_4\}$  and assume that the vectors in subset  $X_t$  are distributed among the classes with equal probability.



# Looking for Homogeneity!!!

## In order for the tree growing methodology

From the root node down to the leaves every split must generate a subsets that are more homogeneous compared to the ancestor's subset  $X_t$ .

## Meaning

The training feature vectors in each one of the new subsets show, whereas data in  $X_t$  are more equally distributed among the classes.

## For example

Consider the task of classifying four classes  $\{\omega_1, \omega_2, \omega_3, \omega_4\}$  and assume that the vectors in subset  $X_t$  are distributed among the classes with equal probability.





# Thus

If we split the node so

- $\omega_1$  and  $\omega_2$  form  $X_{tY}$
- $\omega_3$  and  $\omega_4$  form  $X_{tN}$

Then

$X_{tY}$  and  $X_{tN}$  are more homogeneous compared to  $X_t$ .

in other words

"Purer" in the decision tree terminology.



# Thus

If we split the node so

- $\omega_1$  and  $\omega_2$  form  $X_{tY}$
- $\omega_3$  and  $\omega_4$  form  $X_{tN}$

Then

$X_{tY}$  and  $X_{tN}$  are more homogeneous compared to  $X_t$ .

In other words

"Purer" in the decision tree terminology.



# Thus

If we split the node so

- $\omega_1$  and  $\omega_2$  form  $X_{tY}$
- $\omega_3$  and  $\omega_4$  form  $X_{tN}$

Then

$X_{tY}$  and  $X_{tN}$  are more homogeneous compared to  $X_t$ .

In other words

“Purer” in the decision tree terminology.



# Our Goal

We need

To define a measure that quantifies node impurity.

This

The Overall Impurity of the descendant nodes is optimally decreased with respect to the ancestor node's impurity.



# Our Goal

## We need

To define a measure that quantifies node impurity.

## Thus

The Overall Impurity of the descendant nodes is optimally decreased with respect to the ancestor node's impurity.



# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- **Probabilistic Impurity**
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



# Probabilistic Impurity

Assume the following probability of a vector in  $\mathcal{X}_i$  belongs to class  $i$ :

$$P(\omega_i|t) \text{ for } i = 1, \dots, M \quad (2)$$



# Probabilistic Impurity

Assume the following probability of a vector in  $X_t$  belongs to class  $\omega_i$

$$P(\omega_i|t) \text{ for } i = 1, \dots, M \quad (2)$$





# A Common Impurity

We define one of the most common impurities

$$I(t) = -\sum_{i=1}^M P(\omega_i|t) \log_2 P(\omega_i|t)$$

This is nothing more than the Shannon's Entropy!!

- Facts:

- ▶  $I(t)$  reaches its maximum when

$$P(\omega_i|t) = \frac{1}{M}$$

# A Common Impurity

We define one of the most common impurities

$$I(t) = -\sum_{i=1}^M P(\omega_i|t) \log_2 P(\omega_i|t)$$

This is nothing more than the Shannon's Entropy!!!

- Facts:

- ▶  $I(t)$  reaches its maximum when

$$P(\omega_i|t) = \frac{1}{M}$$

- ▶  $I(t) = 0$  if all data belongs to a single class i.e.

$P(\omega_i|t) = 1$  for only one class, and  $P(\omega_j|t, j \neq i) = 0$  for everybody else.

# A Common Impurity

We define one of the most common impurities

$$I(t) = -\sum_{i=1}^M P(\omega_i|t) \log_2 P(\omega_i|t)$$

This is nothing more than the Shannon's Entropy!!!

- Facts:

- ▶  $I(t)$  reaches its maximum when

$$P(\omega_i|t) = \frac{1}{M}$$

- ▶  $I(t) = 0$  if all data belongs to a single class i.e.

$P(\omega_i|t) = 1$  for only one class, and  $P(\omega_j|t, j \neq i) = 0$  for everybody else.

In reality...

We estimate

$$P(\omega_i|t) = \frac{N_t^i}{N_t}$$

Where  $|\omega_i| = N_t^i$  as the number of points in  $X_t$  that belongs to class  $\omega_i$ .

Assume now

If we perform a split,  $N_{t,Y}$  points are sent into the "YES" node  $X_{t,Y}$  and  $N_{t,N}$  into the "NO" node  $X_{t,N}$



In reality...

We estimate

$$P(\omega_i|t) = \frac{N_t^i}{N_t}$$

Where  $|\omega_i| = N_t^i$  as the number of points in  $X_t$  that belongs to class  $\omega_i$ .

Assume now

If we perform a split,  $N_{tY}$  points are sent into the "YES" node  $X_{tY}$  and  $N_{tN}$  into the "NO" node  $X_{tN}$



## Decrease in node impurity

Then

In a recursive way we define the term decrease in node impurity as:

$$\Delta I(t) = I(t) - \frac{N_{tY}}{N_t} I(t_Y) - \frac{N_{tN}}{N_t} I(t_N) \quad (3)$$

where  $I(t_Y)$  and  $I(t_N)$  are the impurities of the  $t_Y$  and  $t_N$  nodes.



# The Final Goal

## The Final Goal

To adopt from the set of candidate questions the one that performs the split with the highest decrease of impurity.



# Stop-Splitting Rule

## Now

The natural question that now arises is when one decides to stop splitting a node and declares it as a leaf of the tree.





# Stop-Splitting Rule

## Now

The natural question that now arises is when one decides to stop splitting a node and declares it as a leaf of the tree.

## For example you can adopt

A threshold  $T$  and stop splitting if the maximum value of  $\Delta I(t)$  over all possible splits is less than  $T$ .



# Stop-Splitting Rule

## Now

The natural question that now arises is when one decides to stop splitting a node and declares it as a leaf of the tree.

## For example you can adopt

A threshold  $T$  and stop splitting if the maximum value of  $\Delta I(t)$  over all possible splits is less than  $T$ .

## Other possibilities

- If the subset  $X_t$  is small enough.
- If the subset  $X_t$  is pure, in the sense that all points in it belong to a single class.



# Stop-Splitting Rule

## Now

The natural question that now arises is when one decides to stop splitting a node and declares it as a leaf of the tree.

## For example you can adopt

A threshold  $T$  and stop splitting if the maximum value of  $\Delta I(t)$  over all possible splits is less than  $T$ .

## Other possibilities

- If the subset  $X_t$  is small enough.
- If the subset  $X_t$  is pure, in the sense that all points in it belong to a single class.



Once a node is declared to be a leaf

### Class Assignment Rule

Once a node is declared a leaf, we assign the leaf to a class using the rule:

$$j = \arg \max_i P(\omega_i | t).$$



# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- **Final Algorithm**

## 5 Conclusions

- First Some Remarks
- Issues



# Final Algorithm

## Algorithm

- 1 Begin with the root node, that is,  $X_t = X$ .
- 2 For each new node  $t$ 
  - 3 For every feature  $x_k, k = 1, 2, \dots, l$ :
    - 4 For every value  $\alpha_{kn}, n = 1, 2, \dots, N_{tk}$ 
      - 5 Generate  $X_{tY}$  and  $X_{tN}$  according to the answer in the question:  
"Is  $x_k(i) \leq \alpha_{kn}$ ,"  $i = 1, 2, \dots, N_t$
      - 6 Compute the impurity decrease
      - 7 Choose  $\alpha_{kn_0}$  leading to the maximum decrease w.r. to  $x_k$ .
    - 8 Choose  $x_{k_0}$  and associated  $\alpha_{k_0n_0}$  for overall maximum decrease of impurity.
  - 9 If the stop-splitting rule is met, declare node  $t$  as a leaf and label a class
  - 10 If not, generate two descendant nodes  $t_Y$  and  $t_N$  with subsets  $X_{t_Y}$  and  $X_{t_N}$ 
    - 11 depending on the answer to the question: is  $x_{k_0} \leq \alpha$ ?

# Final Algorithm

## Algorithm

- 1 Begin with the root node, that is,  $X_t = X$ .
- 2 For each new node  $t$ 
  - 3 For every feature  $x_k, k = 1, 2, \dots, l$ :
    - 4 For every value  $\alpha_{kn}, n = 1, 2, \dots, N_{tk}$ 
      - 5 Generate  $X_{tY}$  and  $X_{tN}$  according to the answer in the question:  
"Is  $x_k(i) \leq \alpha_{kn}$ ,"  $i = 1, 2, \dots, N_t$
      - 6 Compute the impurity decrease
      - 7 Choose  $\alpha_{kn_0}$  leading to the maximum decrease w.r. to  $x_k$ .
    - 8 Choose  $x_{k_0}$  and associated  $\alpha_{k_0n_0}$  for overall maximum decrease of impurity.
  - 9 If the stop-splitting rule is met, declare node  $t$  as a leaf and label a class
  - 10 If not, generate two descendant nodes  $t_Y$  and  $t_N$  with subsets  $X_{t_Y}$  and  $X_{t_N}$ 
    - 11 depending on the answer to the question: is  $x_{k_0} \leq \alpha$ ?

# Final Algorithm

## Algorithm

- 1 Begin with the root node, that is,  $X_t = X$ .
- 2 For each new node  $t$
- 3 For every feature  $x_k$ ,  $k = 1, 2, \dots, l$ :
  - 4 For every value  $\alpha_{kn}$ ,  $n = 1, 2, \dots, N_{tk}$
  - 5 Generate  $X_{tY}$  and  $X_{tN}$  according to the answer in the question: "Is  $x_k(i) \leq \alpha_{kn}$ ,"  $i = 1, 2, \dots, N_t$
  - 6 Compute the impurity decrease
  - 7 Choose  $\alpha_{kn_0}$  leading to the maximum decrease w.r. to  $x_k$ .
  - 8 Choose  $x_{k_0}$  and associated  $\alpha_{k_0n_0}$  for overall maximum decrease of impurity.
  - 9 If the stop-splitting rule is met, declare node  $t$  as a leaf and label a class
  - 10 If not, generate two descendant nodes  $t_Y$  and  $t_N$  with subsets  $X_{t_Y}$  and  $X_{t_N}$
  - 11 depending on the answer to the question: is  $x_{k_0} \leq \alpha$ ?



# Final Algorithm

## Algorithm

- 1 Begin with the root node, that is,  $X_t = X$ .
- 2 For each new node  $t$
- 3     For every feature  $x_k$ ,  $k = 1, 2, \dots, l$ :
- 4         For every value  $\alpha_{kn}$ ,  $n = 1, 2, \dots, N_{tk}$ 
  - 5             Generate  $X_{tY}$  and  $X_{tN}$  according to the answer in the question: "Is  $x_k(i) \leq \alpha_{kn}$ ,"  $i = 1, 2, \dots, N_t$
  - 6             Compute the impurity decrease
  - 7             Choose  $\alpha_{kn_0}$  leading to the maximum decrease w.r. to  $x_k$ .
  - 8             Choose  $x_{k_0}$  and associated  $\alpha_{x_{k_0}n_0}$  for overall maximum decrease of impurity.
  - 9             If the stop-splitting rule is met, declare node  $t$  as a leaf and label a class
  - 10            If not, generate two descendant nodes  $t_Y$  and  $t_N$  with subsets  $X_{t_Y}$  and  $X_{t_N}$
  - 11            depending on the answer to the question: is  $x_{k_0} \leq \alpha$ ?

# Final Algorithm

## Algorithm

- 1 Begin with the root node, that is,  $X_t = X$ .
- 2 For each new node  $t$
- 3     For every feature  $x_k$ ,  $k = 1, 2, \dots, l$ :
- 4         For every value  $\alpha_{kn}$ ,  $n = 1, 2, \dots, N_{tk}$
- 5             Generate  $X_{tY}$  and  $X_{tN}$  according to the answer in the question:  
              "Is  $x_k(i) \leq \alpha_{kn}$ ,"  $i = 1, 2, \dots, N_t$   
              Compute the impurity decrease  
              Choose  $\alpha_{kN_0}$  leading to the maximum decrease w.r. to  $x_k$ .  
              Choose  $x_{k_0}$  and associated  $\alpha_{x_{k_0}N_0}$  for overall maximum decrease of impurity.  
              If the stop-splitting rule is met, declare node  $t$  as a leaf and label a class  
              If not, generate two descendant nodes  $t_Y$  and  $t_N$  with subsets  $X_{tY}$  and  $X_{tN}$   
              depending on the answer to the question: is  $x_{k_0} \leq \alpha$ ?

# Final Algorithm

## Algorithm

- 1 Begin with the root node, that is,  $X_t = X$ .
- 2 For each new node  $t$
- 3     For every feature  $x_k$ ,  $k = 1, 2, \dots, l$ :
- 4         For every value  $\alpha_{kn}$ ,  $n = 1, 2, \dots, N_{tk}$
- 5             Generate  $X_{tY}$  and  $X_{tN}$  according to the answer in the question:  
6                 "Is  $x_k(i) \leq \alpha_{kn}$ ,"  $i = 1, 2, \dots, N_t$
- 7                 Compute the impurity decrease
- 8                 Choose  $\alpha_{kN_0}$  leading to the maximum decrease w.r. to  $x_k$ .
- 9             Choose  $x_{k_0}$  and associated  $\alpha_{k_0N_0}$  for overall maximum decrease of impurity.
- 10            If the stop-splitting rule is met, declare node  $t$  as a leaf and label a class
- 11            If not, generate two descendant nodes  $t_Y$  and  $t_N$  with subsets  $X_{t_Y}$  and  $X_{t_N}$
- 12            depending on the answer to the question: is  $x_{k_0} \leq \alpha$ ?

# Final Algorithm

## Algorithm

- 1 Begin with the root node, that is,  $X_t = X$ .
- 2 For each new node  $t$
- 3     For every feature  $x_k$ ,  $k = 1, 2, \dots, l$ :
- 4         For every value  $\alpha_{kn}$ ,  $n = 1, 2, \dots, N_{tk}$
- 5             Generate  $X_{tY}$  and  $X_{tN}$  according to the answer in the question:  
6                 "Is  $x_k(i) \leq \alpha_{kn}$ ,"  $i = 1, 2, \dots, N_t$
- 7             Compute the impurity decrease
- 8             Choose  $\alpha_{kN_0}$  leading to the maximum decrease w.r. to  $x_k$ .
- 9             Choose  $x_{k_0}$  and associated  $\alpha_{k_0N_0}$  for overall maximum decrease of impurity.
- 10            If the stop-splitting rule is met, declare node  $t$  as a leaf and label a class
- 11            If not, generate two descendant nodes  $t_Y$  and  $t_N$  with subsets  $X_{t_Y}$  and  $X_{t_N}$
- 12            depending on the answer to the question: is  $x_{k_0} \leq \alpha$ ?

# Final Algorithm

## Algorithm

- 1 Begin with the root node, that is,  $X_t = X$ .
- 2 For each new node  $t$
- 3     For every feature  $x_k$ ,  $k = 1, 2, \dots, l$ :
- 4         For every value  $\alpha_{kn}$ ,  $n = 1, 2, \dots, N_{tk}$
- 5             Generate  $X_{tY}$  and  $X_{tN}$  according to the answer in the question:  
6                 "Is  $x_k(i) \leq \alpha_{kn}$ ,"  $i = 1, 2, \dots, N_t$
- 7             Compute the impurity decrease
- 8             Choose  $\alpha_{kn_0}$  leading to the maximum decrease w.r. to  $x_k$ .
- 9             Choose  $x_{k_0}$  and associated  $\alpha_{k_0 n_0}$  for overall maximum decrease of impurity.
- 10            If the stop-splitting rule is met, declare node  $t$  as a leaf and label a class
- 11            If not, generate two descendant nodes  $t_Y$  and  $t_N$  with subsets  $X_{t_Y}$  and  $X_{t_N}$
- 12            depending on the answer to the question: is  $x_{k_0} \leq \alpha$ ?

# Final Algorithm

## Algorithm

- 1 Begin with the root node, that is,  $X_t = X$ .
- 2 For each new node  $t$
- 3     For every feature  $x_k$ ,  $k = 1, 2, \dots, l$ :
- 4         For every value  $\alpha_{kn}$ ,  $n = 1, 2, \dots, N_{tk}$
- 5             Generate  $X_{tY}$  and  $X_{tN}$  according to the answer in the question:  
6                 "Is  $x_k(i) \leq \alpha_{kn}$ ,"  $i = 1, 2, \dots, N_t$
- 7             Compute the impurity decrease
- 8             Choose  $\alpha_{kn_0}$  leading to the maximum decrease w.r. to  $x_k$ .
- 9     Choose  $x_{k_0}$  and associated  $\alpha_{k_0 n_0}$  for overall maximum decrease of impurity.
- 10     If the stop-splitting rule is met, declare node  $t$  as a leaf and label a class
- 11     If not, generate two descendant nodes  $t_Y$  and  $t_N$  with subsets  $X_{t_Y}$  and  $X_{t_N}$
- 12     depending on the answer to the question: is  $x_{k_0} \leq \alpha$ ?

# Final Algorithm

## Algorithm

- 1 Begin with the root node, that is,  $X_t = X$ .
- 2 For each new node  $t$
- 3     For every feature  $x_k$ ,  $k = 1, 2, \dots, l$ :
- 4         For every value  $\alpha_{kn}$ ,  $n = 1, 2, \dots, N_{tk}$
- 5             Generate  $X_{tY}$  and  $X_{tN}$  according to the answer in the question:  
6                 "Is  $x_k(i) \leq \alpha_{kn}$ ,"  $i = 1, 2, \dots, N_t$
- 7             Compute the impurity decrease
- 8             Choose  $\alpha_{kn_0}$  leading to the maximum decrease w.r. to  $x_k$ .
- 9     Choose  $x_{k_0}$  and associated  $\alpha_{k_0 n_0}$  for overall maximum decrease of impurity.
- 10  If the stop-splitting rule is met, declare node  $t$  as a leaf and label a class
- If not, generate two descendant nodes  $t_Y$  and  $t_N$  with subsets  $X_{t_Y}$  and  $X_{t_N}$
- depending on the answer to the question: is  $x_{k_0} \leq \alpha$ ?

# Final Algorithm

## Algorithm

- 1 Begin with the root node, that is,  $X_t = X$ .
- 2 For each new node  $t$
- 3     For every feature  $x_k$ ,  $k = 1, 2, \dots, l$ :
- 4         For every value  $\alpha_{kn}$ ,  $n = 1, 2, \dots, N_{tk}$
- 5             Generate  $X_{tY}$  and  $X_{tN}$  according to the answer in the question:  
6                 "Is  $x_k(i) \leq \alpha_{kn}$ ,"  $i = 1, 2, \dots, N_t$
- 7             Compute the impurity decrease
- 8             Choose  $\alpha_{kn_0}$  leading to the maximum decrease w.r. to  $x_k$ .
- 9     Choose  $x_{k_0}$  and associated  $\alpha_{k_0 n_0}$  for overall maximum decrease of impurity.
- 10    If the stop-splitting rule is met, declare node  $t$  as a leaf and label a class
- 11    If not, generate two descendant nodes  $t_Y$  and  $t_N$  with subsets  $X_{tY}$  and  $X_{tN}$
- 12         depending on the answer to the question: is  $x_{k_0} \leq \alpha$ ?



# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- Issues



## Popular Classification Methods

- Decision trees have emerged as one of the most popular methods of classification.

### More Impurity Measures

- A variety of node impurity measures can be defined.

### The size of the trees need to be controlled

- The threshold  $T$  leads incorrect sizes.



## Popular Classification Methods

- Decision trees have emerged as one of the most popular methods of classification.

## More Impurity Measures

- A variety of node impurity measures can be defined.

The size of the tree needs to be controlled.

- The threshold  $T$  leads incorrect sizes.



## Popular Classification Methods

- Decision trees have emerged as one of the most popular methods of classification.

## More Impurity Measures

- A variety of node impurity measures can be defined.

## The size of the three need to be controlled

- The threshold  $T$  leads incorrect sizes.



## Why Binary Splits?

- We could consider a Multi-way split

However

- That will fragment the data too fast.

We would rather do only split when necessary

- After all a Multi-way split can be achieved with multiple binary split.



## Why Binary Splits?

- We could consider a Multi-way split

## However

- That will fragment the data too fast.

We would rather do only split when necessary.

- After all a Multi-way split can be achieved with multiple binary split.



## Why Binary Splits?

- We could consider a Multi-way split

## However

- That will fragment the data too fast.

## We would rather do only split when necessary

- After all a Multi-way split can be achieved with multiple binary split.



## Remark

### Linear Combination Splits

- Instead of doing simple splittings, we could use

$$\sum_{j=1}^d a_j x_j < s$$

This improve the predictive power of the tree

- It can hurts interpretability

Similar use

- Hierarchical Mixture of Experts (HME).





## Remark

### Linear Combination Splits

- Instead of doing simple splittings, we could use

$$\sum_{j=1}^d a_j x_j < s$$

### This improve the predictive power of the tree

- It can hurts interpretability

### Summary

- Hierarchical Mixture of Experts (HME).



## Remark

### Linear Combination Splits

- Instead of doing simple splittings, we could use

$$\sum_{j=1}^d a_j x_j < s$$

### This improve the predictive power of the tree

- It can hurts interpretability

### Better use

- Hierarchical Mixture of Experts (HME).



# Outline

## 1 First Principles, Marcus Aurelius (Circa 170 AD)

- Introduction
- Examples of Trees

## 2 Decision Trees

- Deriving Why do they work?
- Structure of Decision Trees
- Types of Decision Trees

## 3 Regression Trees

- Growing Regression Trees
- Using the Sum of Squared Error
- Pruning

## 4 Classification Trees

- Definition
- Training
- The Sought Criterion
- Probabilistic Impurity
- Final Algorithm

## 5 Conclusions

- First Some Remarks
- **Issues**



## One of the biggest issues

- One major problem with trees is their high variance.

A small change in the data can result in a very different series of splits

- Making interpretability precarious!!!



# Issues

## One of the biggest issues

- One major problem with trees is their high variance.

## A small change in the data can result in a very different series of splits

- Making interpretability precarious!!!



## Lack of Smoothness

- Another limitation of trees is the lack of smoothness of the prediction surface

• Thus strategies to alleviate this problem are necessary

- Multivariate Adaptive Regression Splines (MARS) procedure



## Lack of Smoothness

- Another limitation of trees is the lack of smoothness of the prediction surface

Thus strategies to alleviate this problem are necessary

- Multivariate Adaptive Regression Splines (MARS) procedure



## The CART trees are bad at modeling additive structures

- For Example

$$y = c_1 I(x_1 < t_1) + c_2 I(x_2 < t_2) + \epsilon \text{ with } \epsilon \sim N(0, \sigma^2)$$

Problem: CART has no special encouragement to capture this model

- Again MARS can help for this given its no dependency to the binary tree structure





## The CART trees are bad at modeling additive structures

- For Example

$$y = c_1 I(x_1 < t_1) + c_2 I(x_2 < t_2) + \epsilon \text{ with } \epsilon \sim N(0, \sigma^2)$$

## Problem, CART has no special encouragement to capture this model

- Again MARS can help for this given its no dependency to the binary tree structure

