Machine Learning for Official Statistics & SDGs

# Decision Trees
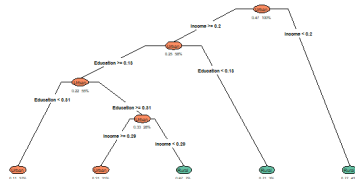
UNITED NATIONS
SIAP
Statistical Institute for
Asia and the Pacific

# [ DECISION TREES ]

*Trees* are method for classification or regression analysis.
↪ Focus on classification



Decision tree with a max depth of 5

## [ DECISION TREES ]

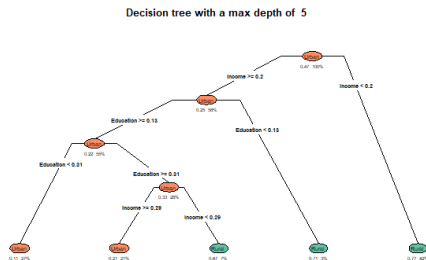*Trees* are method for classification or regression analysis.
$\hookrightarrow$ Focus on classification



Decision tree with a max depth of 5

► Trees split the space into non-overlapping spaces
► Used to assign/predict a *class* following conditions
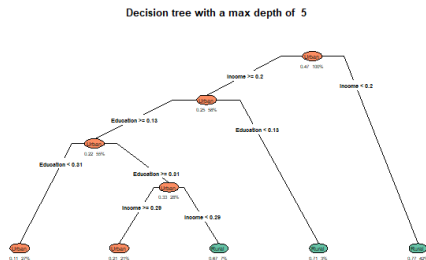► Optimally, final classification is homogenous

[ WHAT'S IN A TREE? ]

Trees have a very simple structure and are easy to understand:



Decision tree with a max depth of 5

## [ WHAT'S IN A TREE? ]

Trees have a very simple structure and are easy to understand:
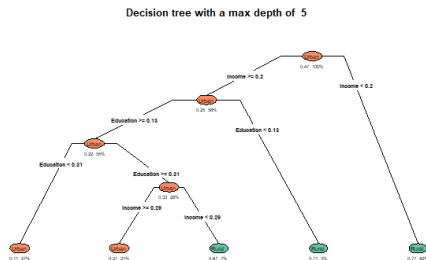


Decision tree with a max depth of 5

Trees have:
▶ **Nodes** where splitting decisions are done

# [ WHAT'S IN A TREE? ]

Trees have a very simple structure and are easy to understand:
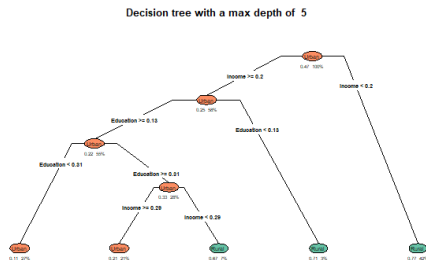


Decision tree with a max depth of 5

Trees have:
- ▶ **Nodes** where splitting decisions are done
- ▶ **Branches** following conditions

# [ WHAT'S IN A TREE? ]

Trees have a very simple structure and are easy to understand:
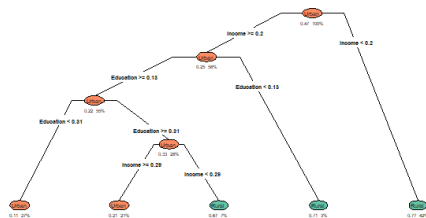


Decision tree with a max depth of 5

Trees have:

- ▶ **Nodes** where splitting decisions are done
- ▶ **Branches** following conditions
- ▶ **Leaves** are terminal nodes of the classification
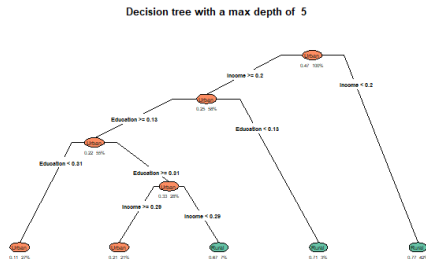
# [ WHAT'S IN A TREE? ]

How to build a tree?



Decision tree with a max depth of 5

# [ WHAT'S IN A TREE? ]

How to build a tree?
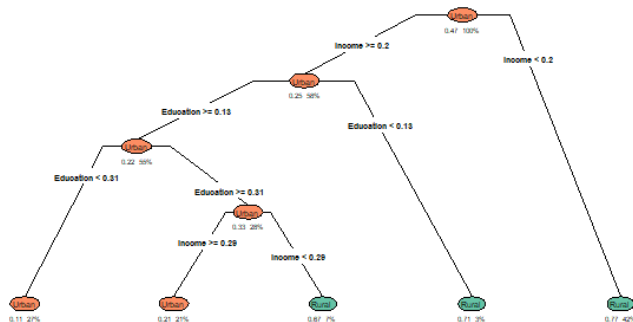


Decision tree with a max depth of 5

- ▶ Trees are based on recursive binary splits
- ▶ Each node uses a threshold on a variable
- ▶ Each node separates the observations in two sets

# [ EXAMPLE ON A SIMPLE TREE ]

Let us see how this tree is constructed:
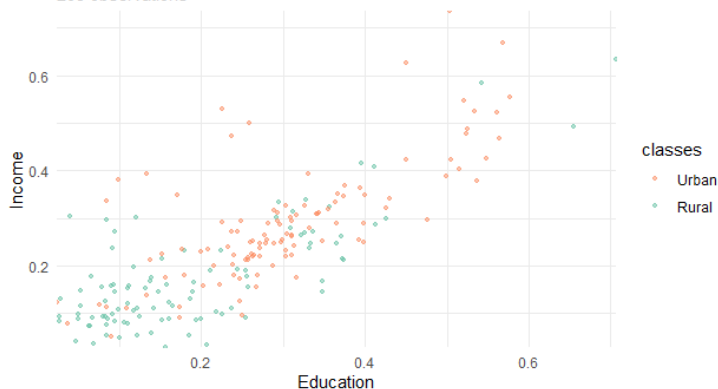


Decision tree with a max depth of 5

4 nodes leading to final leaves ↪ Depth = 5

# [ EXAMPLE ON A SIMPLE TREE ]

## The problem is a 2D space



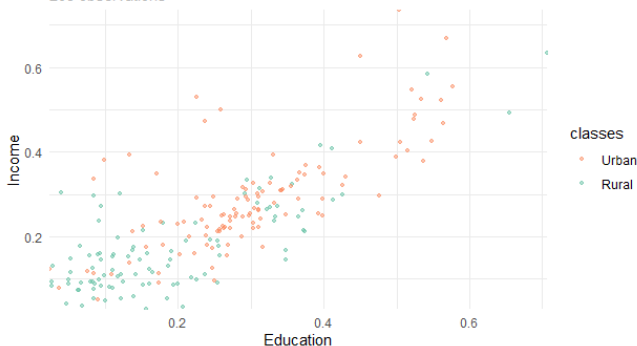Position of Rural and Urban households in (Education, Income) space
208 observations

# [ EXAMPLE ON A SIMPLE TREE ]



Decision tree with a max depth of 5

# [ EXAMPLE ON A SIMPLE TREE ]



Decision tree with a max depth of 5



Position of Rural and Urban households in (Education, Income) space

208 observations
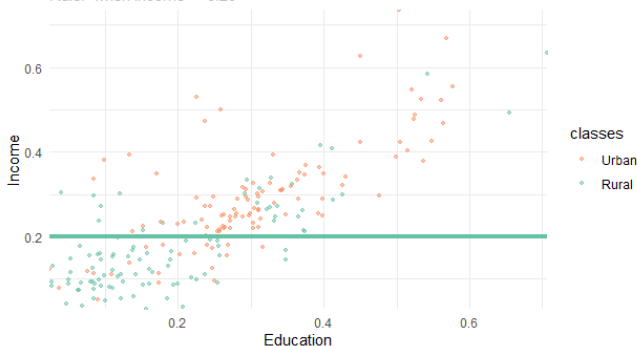
classes
- Urban
- Rural

How to split the (Education,Income) space?

# [ EXAMPLE ON A SIMPLE TREE ]



**First node**
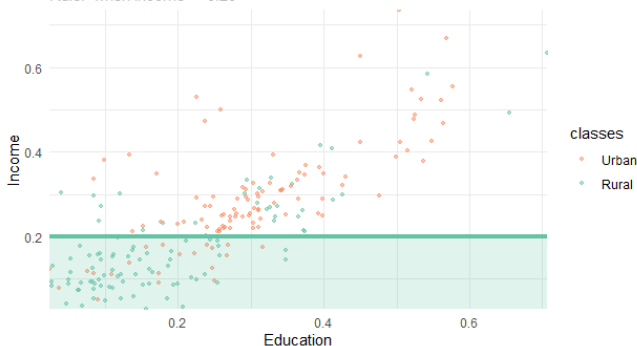Rule: when Income < 0.20



The first boundary decision line

# [ EXAMPLE ON A SIMPLE TREE ]



**First node**

Rule: when Income < 0.20



The space below the line is classified as rural

# [ EXAMPLE ON A SIMPLE TREE ]



Decision tree with a max depth of 5

**Second node:**

Rule: when Income >= 0.20 & Education < 0.13



classes
- Urban
- Rural

Second boundary decision line

# [ EXAMPLE ON A SIMPLE TREE ]



The space on the left of the line is classified as rural

# [ EXAMPLE ON A SIMPLE TREE ]



Third boundary decision line

# [ EXAMPLE ON A SIMPLE TREE ]



**Third node**

Rule: when Income is 0.20 to 0.29 & Education >= 0.31



The space on the right of the line is classified as Urban

# [ EXAMPLE ON A SIMPLE TREE ]



Fourth boundary decision line

# [ EXAMPLE ON A SIMPLE TREE ]



**Fourth node rule**

Rule: when Income >= 0.29 & Education >= 0.31



classes
- Urban
- Rural

The space above the line is classified as Urban

# [ EXAMPLE ON A SIMPLE TREE ]



Finally, the remaining space is classified as rural

# [ HOW TO BUILD A TREE? ]



Decision tree with a max depth of 3

We need several tools to build a tree:

# [ HOW TO BUILD A TREE? ]



Decision tree with a max depth of 3

We need several tools to build a tree:

▶ A method to choose the decision variable (one per node)

# [ HOW TO BUILD A TREE? ]



Decision tree with a max depth of 3

We need several tools to build a tree:

► A method to choose the decision variable (one per node)

► A criterion to define best threshold

[ HOW TO BUILD A TREE? ]



Decision tree with a max depth of 3

We need several tools to build a tree:

- ▶ A method to choose the decision variable (one per node)
- ▶ A criterion to define best threshold
- ▶ A criterion to measure the quality of each split

# [ HOW TO BUILD A TREE? ]



Decision tree with a max depth of 3

We need several tools to build a tree:

►  A method to choose the decision variable (one per node)

►  A criterion to define best threshold

►  A criterion to measure the quality of each split

►  A way to decide when to stop (the terminal node becomes a leaf)

[ TOOLS TO BUILD A TREE ]

At each node, one can measure the *purity* of each split

[ TOOLS TO BUILD A TREE ]

At each node, one can measure the *purity* of each split

▶ Misclassification error rate

## [ TOOLS TO BUILD A TREE ]

At each node, one can measure the *purity* of each split

- ▶ Misclassification error rate
- ▶ The Gini coefficient measures the purity in each node $\kappa$

$$D_\kappa = \sum_{m=1}^{M} \widehat{p}_{m\kappa}(1 - \widehat{p}_{m\kappa})$$

where $\widehat{p}_{m\kappa}$ is the proportion of class $m$ in node $\kappa$.

## [ TOOLS TO BUILD A TREE ]

At each node, one can measure the *purity* of each split

- ▶ Misclassification error rate
- ▶ The Gini coefficient measures the purity in each node $\kappa$

$$D_\kappa = \sum_{m=1}^{M} \widehat{p}_{m\kappa}(1 - \widehat{p}_{m\kappa})$$

where $\widehat{p}_{m\kappa}$ is the proportion of class $m$ in node $\kappa$.

- ▶ The entropy or information:

$$D_\kappa = -\sum_{m=1}^{M} \widehat{p}_{m\kappa} \log \widehat{p}_{m\kappa}$$

## [ TOOLS TO BUILD A TREE ]

At each node, one can measure the *purity* of each split

- ▶ Misclassification error rate
- ▶ The Gini coefficient measures the purity in each node $\kappa$

$$D_\kappa = \sum_{m=1}^{M} \widehat{p}_{m\kappa}(1 - \widehat{p}_{m\kappa})$$

where $\widehat{p}_{m\kappa}$ is the proportion of class $m$ in node $\kappa$.

- ▶ The entropy or information:

$$D_\kappa = -\sum_{m=1}^{M} \widehat{p}_{m\kappa} \log \widehat{p}_{m\kappa}$$

- ▶ We expect that there is an Information gain from the splitting

*Information Gain = Entropy $_{Before}$ − Entropy $_{After}$*

# [ EXAMPLE ON A SIMPLE TREE ]



Decision tree with a max depth of 5

# [ EXAMPLE ON A SIMPLE TREE ]



First the boundary decision line

# [ EXAMPLE ON A SIMPLE TREE ]



The space below the line is classified as rural

# [ EXAMPLE ON A SIMPLE TREE ]



Some "Urban" are classified as rural ↪ Impurity

# [ EXAMPLE ON A SIMPLE TREE ]



Some "Rural" are classified as Urban $\hookrightarrow$ Impurity

# [ EXAMPLE ON A SIMPLE TREE ]



Second boundary decision line

# [ EXAMPLE ON A SIMPLE TREE ]



The space on the left of the line is classified as rural

# [ EXAMPLE ON A SIMPLE TREE ]



**Second node:**

Rule: when Income >= 0.20 & Education < 0.13



Some "Urban" are classified as Rural

# [ EXAMPLE ON A SIMPLE TREE ]



Decision tree with a max depth of 5

**Second node:**
Rule: when Income >= 0.20 & Education < 0.13



classes
· Urban
· Rural

Some "Rural" are classified as Urban

# [ EXAMPLE ON A SIMPLE TREE ]



Third boundary decision line

# [ EXAMPLE ON A SIMPLE TREE ]



**Third node**

Rule: when Income is 0.20 to 0.29 & Education >= 0.31



The space on the right of the line is classified as Urban

# [ EXAMPLE ON A SIMPLE TREE ]



Decision tree with a max depth of 5



**Third node**

Rule: when Income is 0.20 to 0.29 & Education >= 0.31

classes
- Urban
- Rural

Some "Rural" are classified as Urban

# [ EXAMPLE ON A SIMPLE TREE ]

# [ EXAMPLE ON A SIMPLE TREE ]



Fourth boundary decision line

Introduction
○

What's in a tree?
○○

Step-by-Step
○○○

How to build a tree?
○○●○

Tuning a Tree
○○○○○○○○

Wrap-up
○

# [ EXAMPLE ON A SIMPLE TREE ]



Decision tree with a max depth of 5

**Fourth node rule**

Rule: when Income >= 0.29 & Education >= 0.31



classes

• Urban

• Rural

The space above the line is classified as Urban

# [ EXAMPLE ON A SIMPLE TREE ]



Finally, the remaining space is classified as rural

# [ EXAMPLE ON A SIMPLE TREE ]



Few "Rural" are classified as Urban

# [ EXAMPLE ON A SIMPLE TREE ]



Decision tree with a max depth of 5

**Fourth node rule**

Rule: when Income >= 0.29 & Education >= 0.31



classes
- Urban
- Rural

Many "Urban" are classified as Rural

[ HOW TO BUILD A TREE? ]

The construction is based on recursive binary splits

[ HOW TO BUILD A TREE? ]

The construction is based on recursive binary splits

▶ The goal is to increase the quality of the classification at the each stage

## [ HOW TO BUILD A TREE? ]

The construction is based on recursive binary splits

- ▶ The goal is to increase the quality of the classification at the each stage
- ↪ decrease the *impurity* at each node

[ HOW TO BUILD A TREE? ]

The construction is based on recursive binary splits

- ▶ The goal is to increase the quality of the classification at the each stage
- ↪ decrease the *impurity* at each node
- ▶ The outcome is "one" tree, not the perfect one

[ HOW TO BUILD A TREE? ]

The construction is based on recursive binary splits

- ▶ The goal is to increase the quality of the classification at the each stage
- ↪ decrease the *impurity* at each node
- ▶ The outcome is "one" tree, not the perfect one
- ▶ There are many parameters that can be adjusted

# [ HOW TO BUILD A TREE? ]

The construction is based on recursive binary splits

- ▶ The goal is to increase the quality of the classification at the each stage
- ↪ decrease the *impurity* at each node
- ▶ The outcome is "one" tree, not the perfect one
- ▶ There are many parameters that can be adjusted
    - ▶ The maximum *depth* of the tree

# [ HOW TO BUILD A TREE? ]

The construction is based on recursive binary splits

- ▶ The goal is to increase the quality of the classification at the each stage
- ↪ decrease the *impurity* at each node
- ▶ The outcome is "one" tree, not the perfect one
- ▶ There are many parameters that can be adjusted
    - ▶ The maximum *depth* of the tree
    - ▶ The number of final leaves

[ HOW TO BUILD A TREE? ]

The construction is based on recursive binary splits

▶ The goal is to increase the quality of the classification at the each stage

↪ decrease the *impurity* at each node

▶ The outcome is "one" tree, not the perfect one

▶ There are many parameters that can be adjusted

    ▶ The maximum *depth* of the tree
    ▶ The number of final leaves
    ▶ The impurity balance between classes

# [ HOW TO BUILD A TREE? ]

The construction is based on recursive binary splits

- ▶ The goal is to increase the quality of the classification at the each stage
- ↪ decrease the *impurity* at each node
- ▶ The outcome is "one" tree, not the perfect one
- ▶ There are many parameters that can be adjusted
    - ▶ The maximum *depth* of the tree
    - ▶ The number of final leaves
    - ▶ The impurity balance between classes
    - ▶ The *complexity parameter*

[ HOW TO BUILD A TREE? ]

The construction is based on recursive binary splits

▶ The goal is to increase the quality of the classification at the each stage

↪ decrease the *impurity* at each node

▶ The outcome is "one" tree, not the perfect one

▶ There are many parameters that can be adjusted

　▶ The maximum *depth* of the tree
　▶ The number of final leaves
　▶ The impurity balance between classes
　▶ The *complexity parameter*
　▶ ...

# [ TREES CAN BE COMPLEX ]

Trees can decompose the space in very specific zones.
↪ Example with the full set of variables



Decision tree with no constrains

# [ SELECTING THE **DEPTH** OF A TREE ]

Using CV, we can select the *maximum depth* parameter

# [ SELECTING THE **DEPTH** OF A TREE ]

Using CV, we can select the *maximum depth* parameter



Decision Tree - Tuning Maximum Depth

[ SELECTING THE **DEPTH** OF A TREE ]

# [ SELECTING THE **DEPTH** OF A TREE ]

The resulting tree



Tree Selected (max depth = 11)

Tree with optimal depth

# [ SELECTING THE **DEPTH** OF A TREE ]

The resulting tree



Feature importance (also confusion matrix, kappa, etc..)

[ SELECTING THE **COMPLEXITY** OF A TREE ]

The complexity of a tree is a parameter $C_p$ governing the trade-off between tree size $|T|$ and its overall accuracy $D(T)$:

$$D_{C_p}(T) = D(T) + C_p \cdot |T|$$

▶ $D(T) = \sum_{\kappa=1}^{K} D_\kappa$: the total *impurity* of the tree

Introduction    What's in a tree?    Step-by-Step    How to build a tree?    **Tuning a Tree**    Wrap-up

○       ○○       ○○○       ○○○○       ○○○●○○○○       ○

# [ SELECTING THE **COMPLEXITY** OF A TREE ]

The complexity of a tree is a parameter $C_p$ governing the trade-off between tree size $|T|$ and its overall accuracy $D(T)$:

$$D_{C_p}(T) = D(T) + C_p \cdot |T|$$

▶ $D(T) = \sum_{\kappa=1}^{K} D_\kappa$: the total *impurity* of the tree
▶ $|T|$ is the number of terminal nodes of the tree

[ SELECTING THE **COMPLEXITY** OF A TREE ]

The complexity of a tree is a parameter $C_p$ governing the
trade-off between tree size $|T|$ and its overall accuracy $D(T)$:

$$D_{C_p}(T) = D(T) + C_p \cdot |T|$$

▶ $D(T) = \sum_{\kappa=1}^{K} D_\kappa$: the total *impurity* of the tree
▶ $|T|$ is the number of terminal nodes of the tree
↪ A model with $C_p = 0$ will impose no constrains

# [ SELECTING THE **COMPLEXITY** OF A TREE ]

The complexity of a tree is a parameter $C_p$ governing the trade-off between tree size $|T|$ and its overall accuracy $D(T)$:

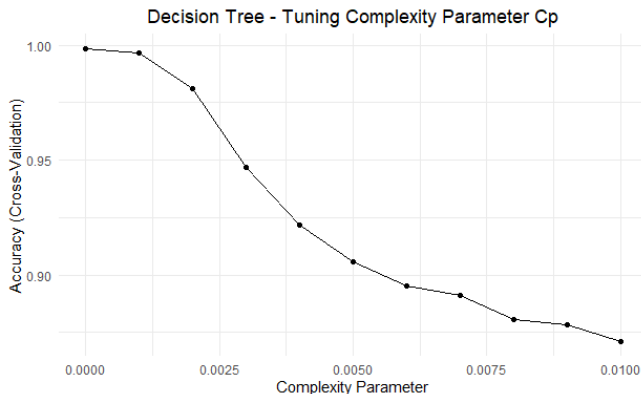$$D_{C_p}(T) = D(T) + C_p \cdot |T|$$

- ► $D(T) = \sum_{\kappa=1}^{K} D_\kappa$: the total *impurity* of the tree
- ► $|T|$ is the number of terminal nodes of the tree
- ↪ A model with $C_p = 0$ will impose no constrains
- ↪ A value of $C_p = 1$ only **one** terminal (and initial) node.

# [ SELECTING THE **COMPLEXITY** OF A TREE ]

# [ SELECTING THE **COMPLEXITY** OF A TREE ]

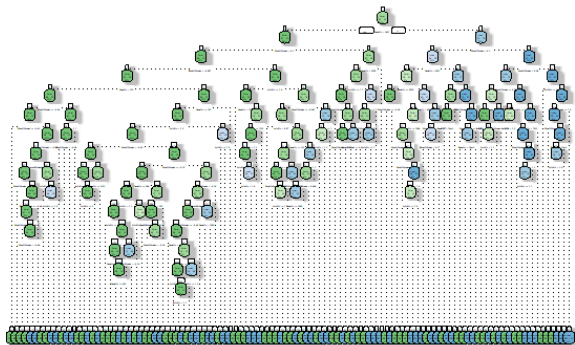The result:



Grid search

# [ SELECTING THE **COMPLEXITY** OF A TREE ]

The result:



Tree with optimized Cp

Final tree with optimal complexity parameter

# [ SELECTING THE **COMPLEXITY** OF A TREE ]

The result:



Feature importance (also confusion matrix, kappa, etc..)

# [ HOW TO **PRUNE** A TREE? ]

One can also *prune* a tree

# [ HOW TO **PRUNE** A TREE? ]

One can also *prune* a tree

1. Final trees may be too large and too complex

[ HOW TO **PRUNE** A TREE? ]

One can also *prune* a tree

1. Final trees may be too large and too complex

↪ Risk of overfitting

[ HOW TO **PRUNE** A TREE? ]

One can also *prune* a tree

1. Final trees may be too large and too complex

↪ Risk of overfitting

2. Pruning techniques use the same criteria on each leave

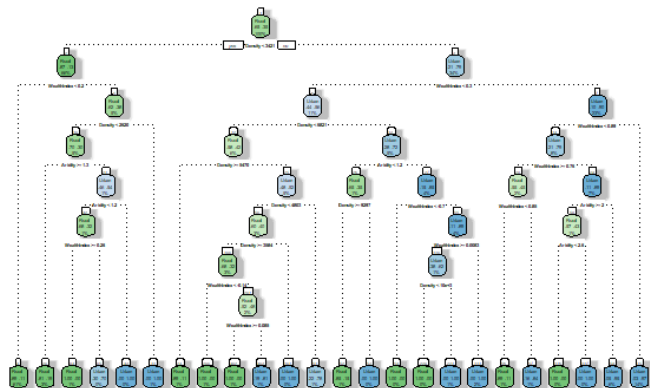# [ HOW TO **PRUNE** A TREE? ]

One can also *prune* a tree

1. Final trees may be too large and too complex

↪ Risk of overfitting

2. Pruning techniques use the same criteria on each leave

↪ Remove least important nodes

# [ PRUNED TREE ]

After "*pruning*":

# [ PRUNED TREE ]

After "*pruning*":



Pruned tree

## [ PRUNED TREE ]

After "*pruning*":



Pruned tree

$\hookrightarrow$ Easier to interpret & no loss of accuracy

# [QUIZ TIME]

# [QUIZ TIME]

# [TAKEAWAYS]

▶ Trees are simple and easy to interpret

# [TAKEAWAYS]

▶ Trees are simple and easy to interpret
▶ Each node is a split based on a threshold

## [TAKEAWAYS]

▶ Trees are simple and easy to interpret

▶ Each node is a split based on a threshold

▶ Splits are determined to maximize some measure of accuracy in each strata of the predictors' space

# [TAKEAWAYS]

- ▶ Trees are simple and easy to interpret
- ▶ Each node is a split based on a threshold
- ▶ Splits are determined to maximize some measure of accuracy in each strata of the predictors' space
- ▶ Regression trees apply the same logic, with different criteria and values

## [TAKEAWAYS]

- ▶ Trees are simple and easy to interpret
- ▶ Each node is a split based on a threshold
- ▶ Splits are determined to maximize some measure of accuracy in each strata of the predictors' space
- ▶ Regression trees apply the same logic, with different criteria and values
- ▶ One can select the depth of a tree or its complexity and prune it

## [TAKEAWAYS]

▶ Trees are simple and easy to interpret

▶ Each node is a split based on a threshold

▶ Splits are determined to maximize some measure of accuracy in each strata of the predictors' space

▶ Regression trees apply the same logic, with different criteria and values

▶ One can select the depth of a tree or its complexity and prune it

▶ Trees are very specific, not robust and prone to overfitting

## [TAKEAWAYS]

▶ Trees are simple and easy to interpret

▶ Each node is a split based on a threshold

▶ Splits are determined to maximize some measure of accuracy in each strata of the predictors' space

▶ Regression trees apply the same logic, with different criteria and values

▶ One can select the depth of a tree or its complexity and prune it

▶ Trees are very specific, not robust and prone to overfitting

↪ There are powerful methods using many trees...