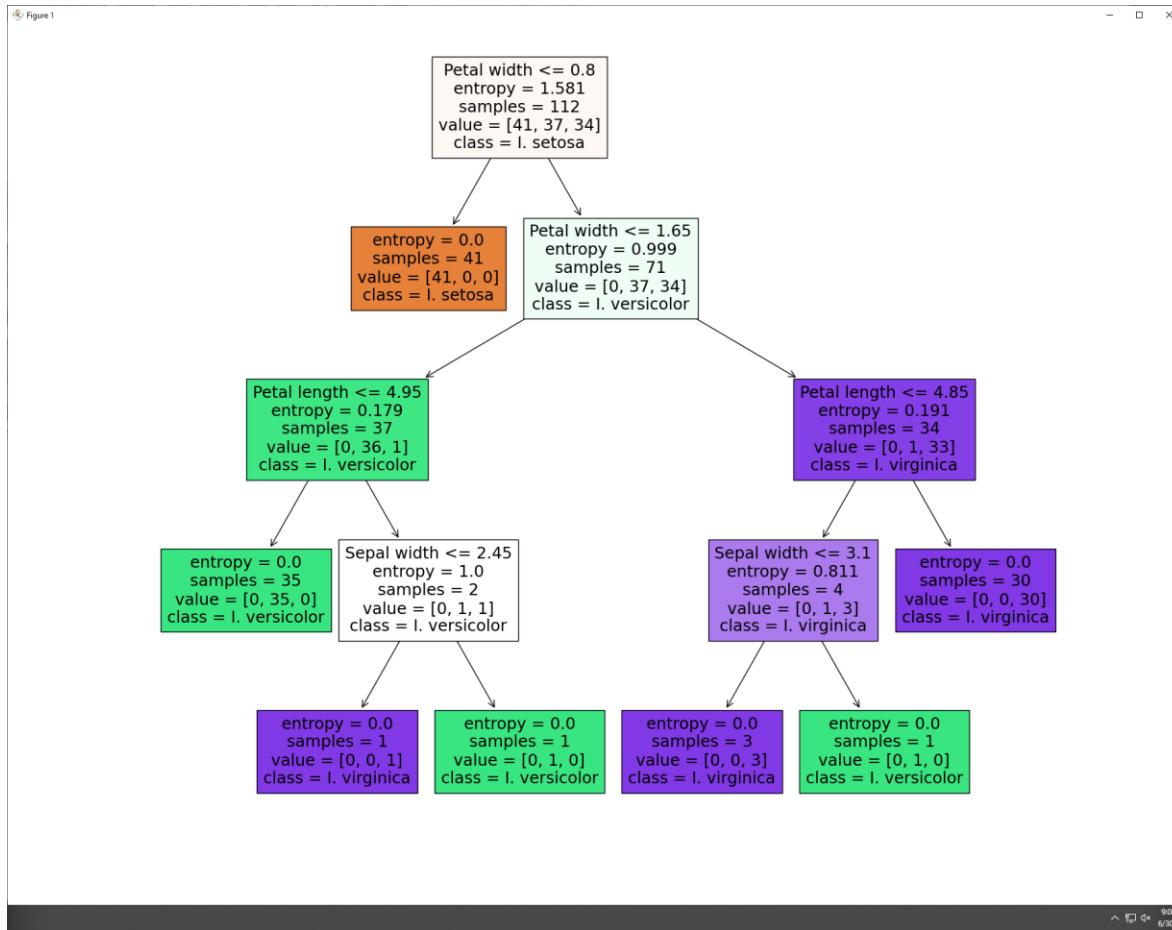


I T P 4 4 9

# Decision Trees

L e c t u r e 1 1





# CART

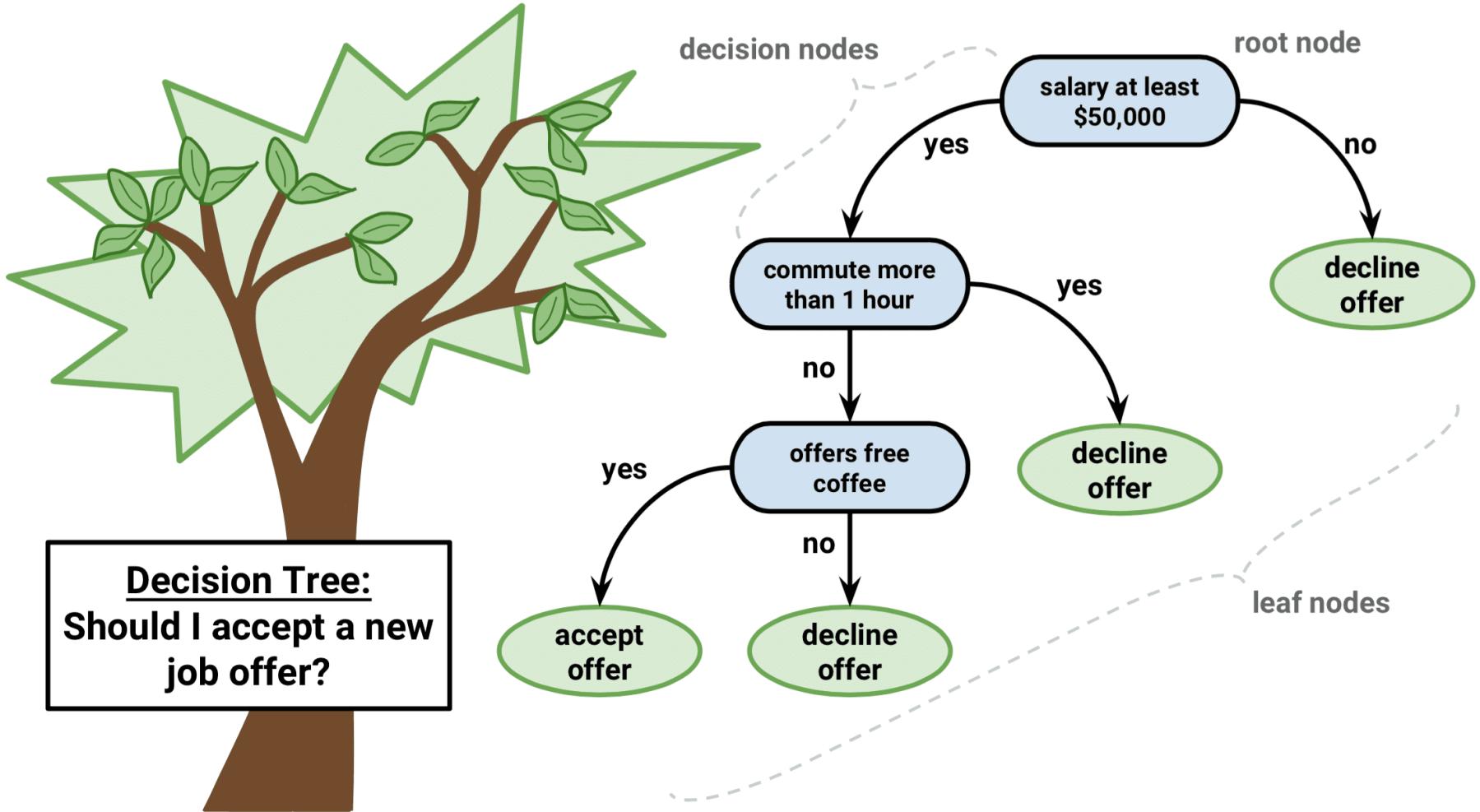
Classification and Regression Trees

# Difference: Classification and Regression

- In **classification** we classify cases into specific **classes**. For example, we might take a set of people and classify them as potential customers or non-customers (based on some predictor attributes)
- In **regression**, we are trying to predict a **numerical** value – like the annual income of a person, the price of a car, the sales volume or GPA or whatever.
- As the name implies, CART can be used for both purposes!
- In this lecture we will only look at the use of CART for **classification**.

# Decision Trees

- *Supervised* classification algorithm that is used when the target variable is a categorical variable and the predictor variables are either categorical or numerical
- Utilizes a set of *if-then* rules, represented as a tree, for classification
- *Visual* output of the model makes it easier to understand and implement compared to other predictive models



# Homogeneity

- The main goal in a decision tree algorithm is to identify a *variable* and *classification* that results in a more homogeneous distribution with reference to the target variable
- The homogeneous distribution means that similar values of the target variable are grouped together so that a *concrete* decision can be made

Income	Education	Family size	Type of car
200,000	1	4	Luxury
30,000	1	4	Compact
75,000	2	3	High-end
90,000	3	6	Economy
30,000	2	2	Compact
125,000	2	2	Luxury
130,000	3	6	Economy
300,000	3	2	Luxury
250,000	2	3	Luxury
128,000	2	5	Economy
35,000	3	1	Compact
...	...	...	...

# Rule Archeology

Income  $\geq 200,000 \rightarrow$  Luxury

Income  $\leq 35,000 \rightarrow$  Compact

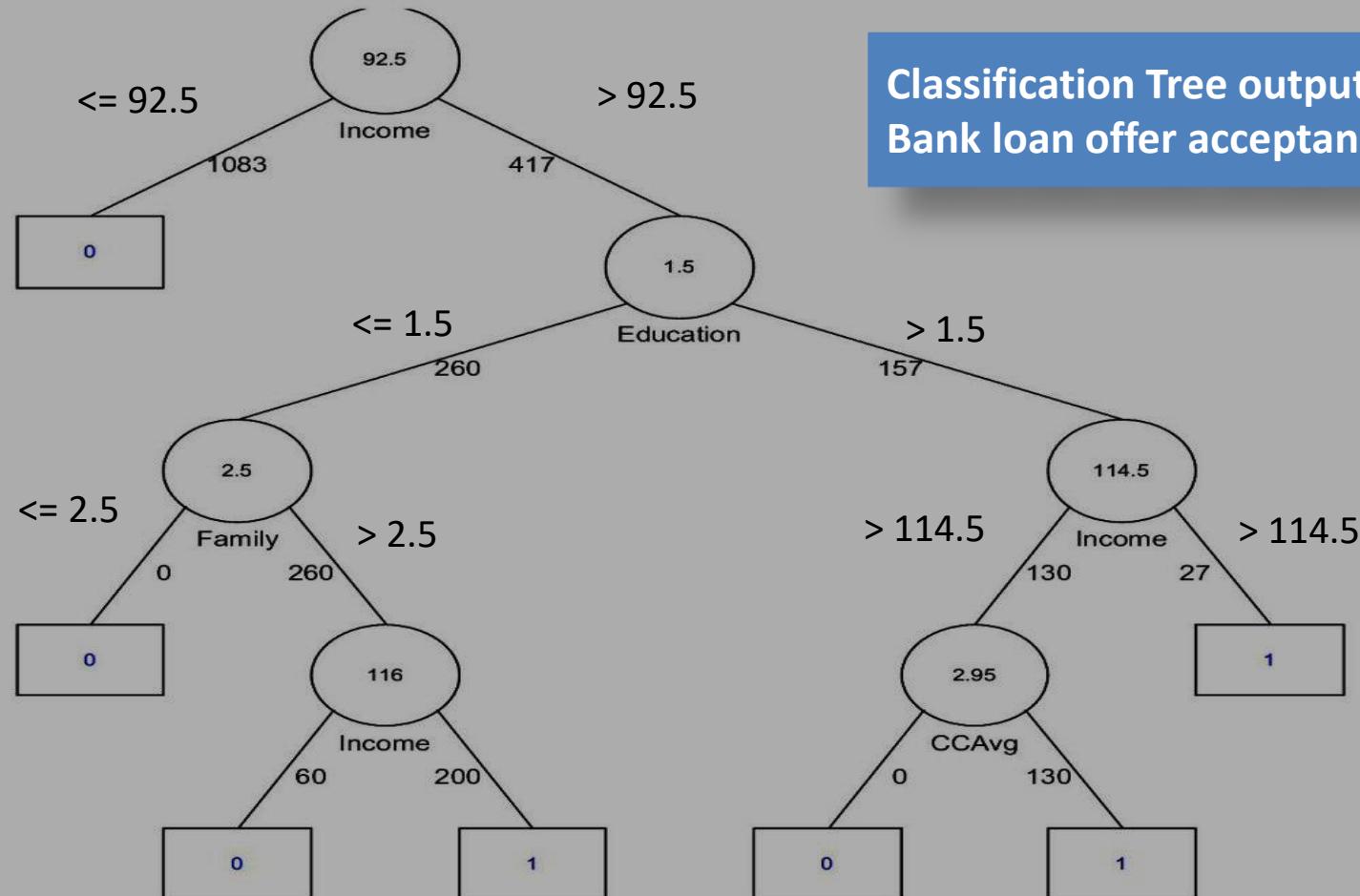
( $100,000 \leq \text{Income} \leq 200,000$ )  
and family size  $< 3 \rightarrow$  Luxury

( $100,000 \leq \text{Income} \leq 200,000$ )  
and family size  $\geq 5 \rightarrow$  Economy

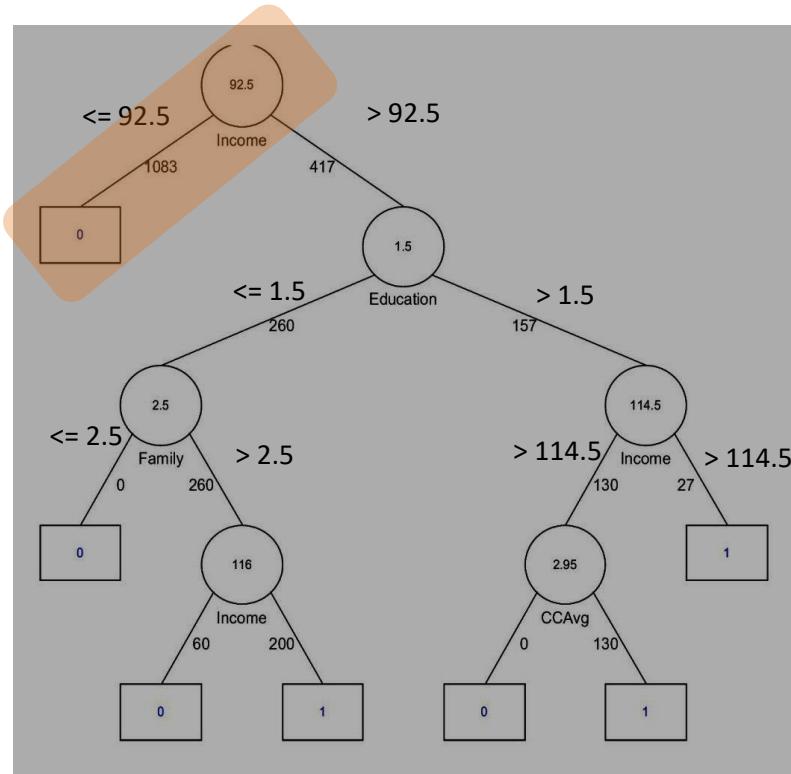
- Looking at the data on the type of car that people own along with other demographic information (made up for illustration), we might be able to infer some underlying “rules” about the kinds of cars that people own.
- With the caveat that we have very little data above and cannot really make any serious inferences, the discussion that follows is only intended to illustrate what we will be studying this week. To really make serious inferences we need much more data.
- From the data we might observe some patterns:
  - People making \$200,000 or more seem to own a luxury car independent of other factors.
  - People making \$35,000 or less seem to be owning a compact car, independent of other factors.
  - For people making between \$100,000 and \$200,000, it seems that they own luxury cars when the family size is 2 or less and own an economy car when the family size is 5 or more.

# **LOAN EXAMPLE**

## Classification Tree output for Bank loan offer acceptance data

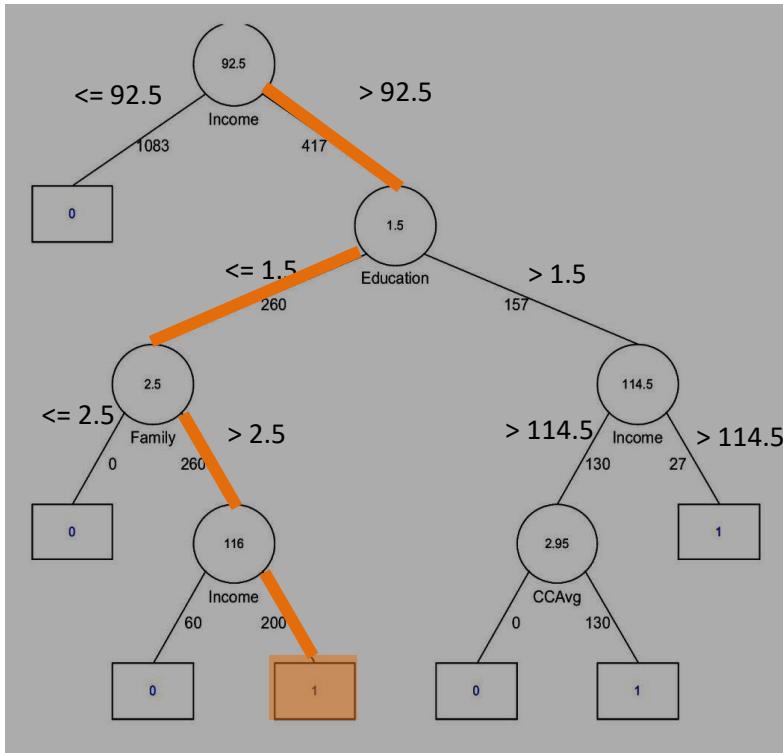


Reference: Data mining for business intelligence: Shmueli, Patel and Bruce, Wiley, 2006



The tree implies rules. The simplest rule is highlighted.

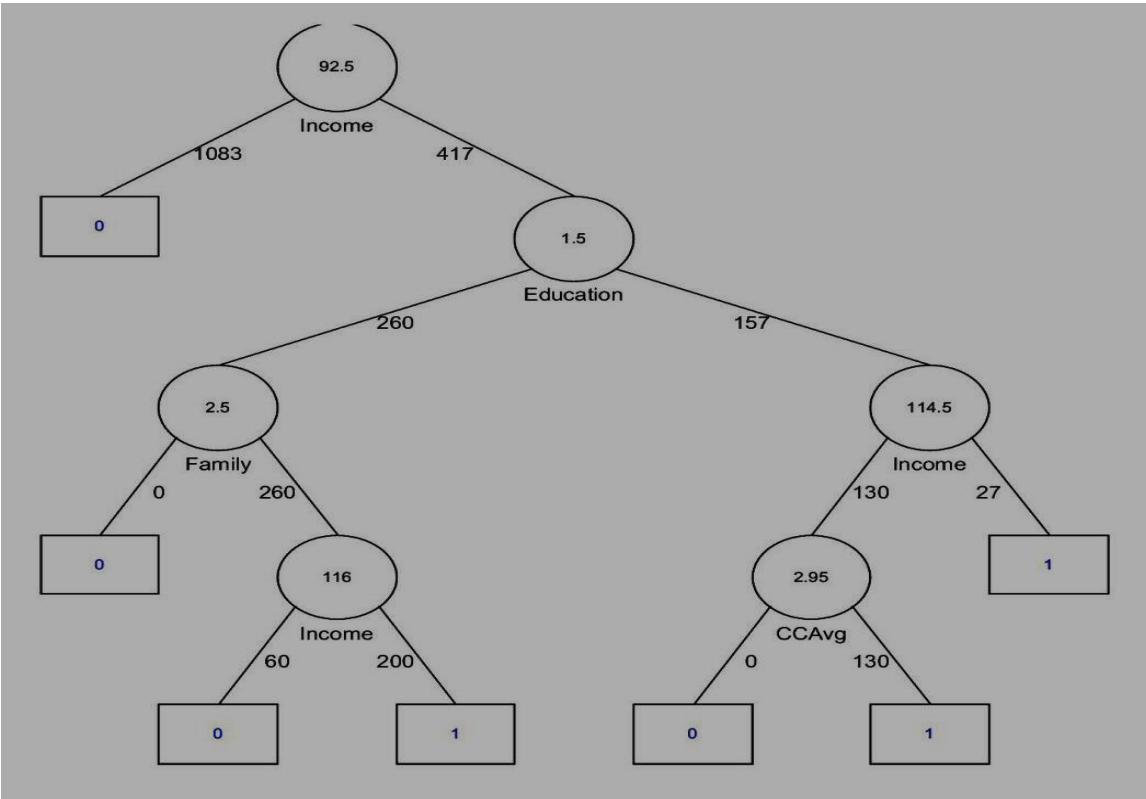
If income  $\leq 92.5$  then Non acceptor



A more complex rule that the tree shows.

If income  $> 92.5$  and  
education  $\leq 1.5$  and  
family size  $> 2.5$  and  
income  $> 114.5$  then  
Acceptor

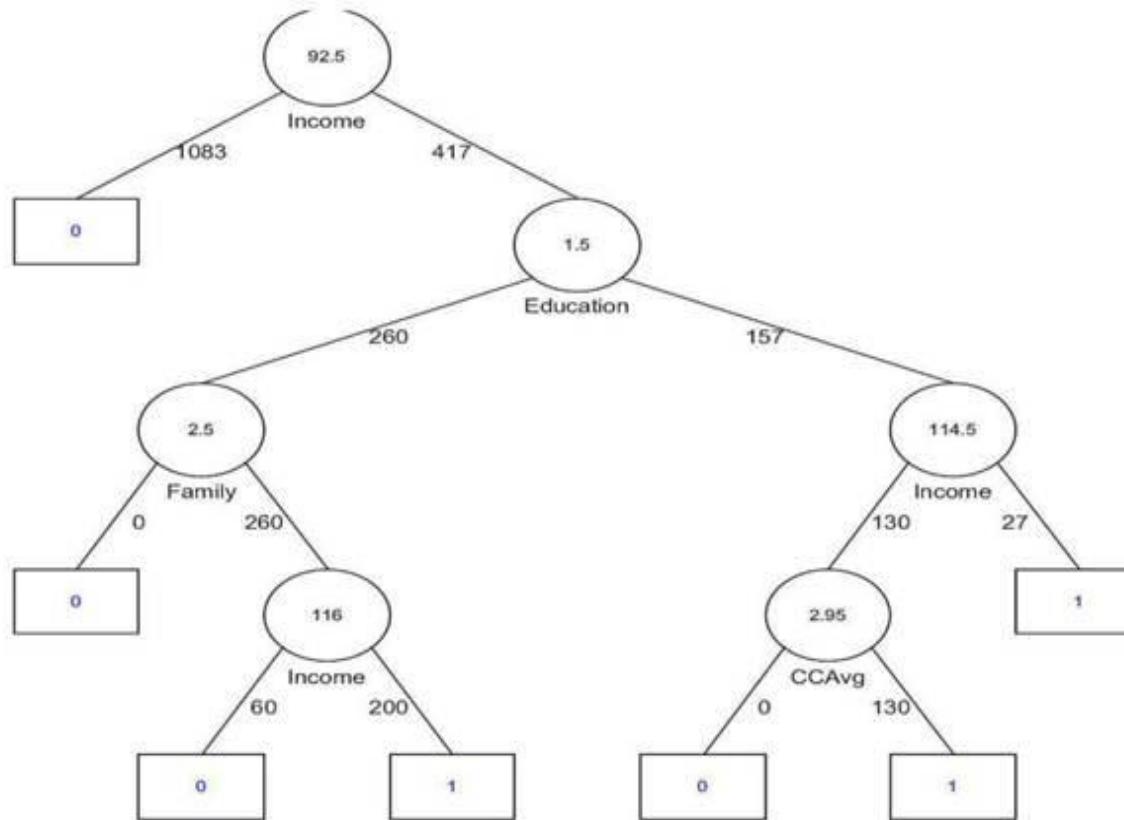
Can you condense  
this rule?

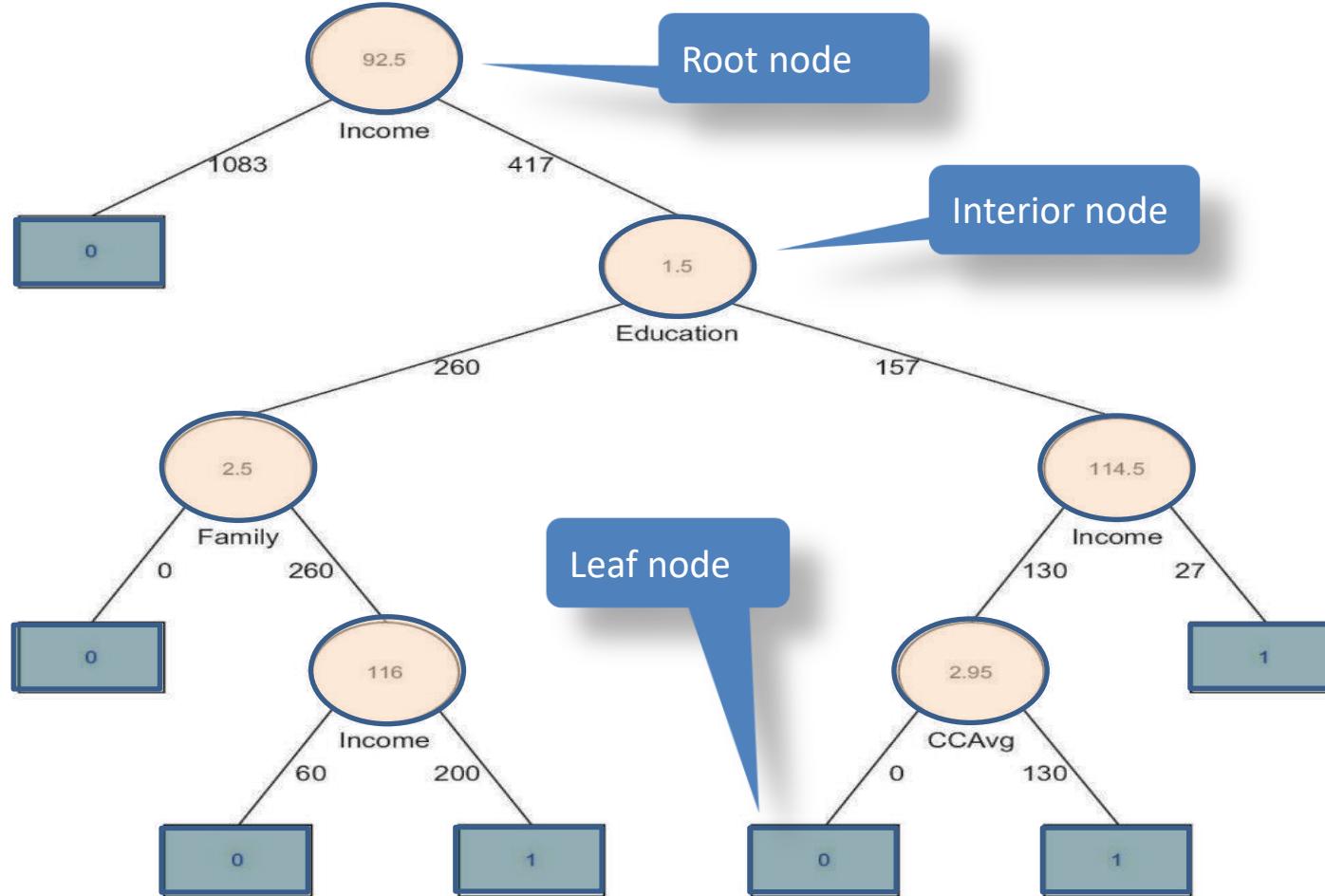


**Find two more rules.**

**How many rules in all?**

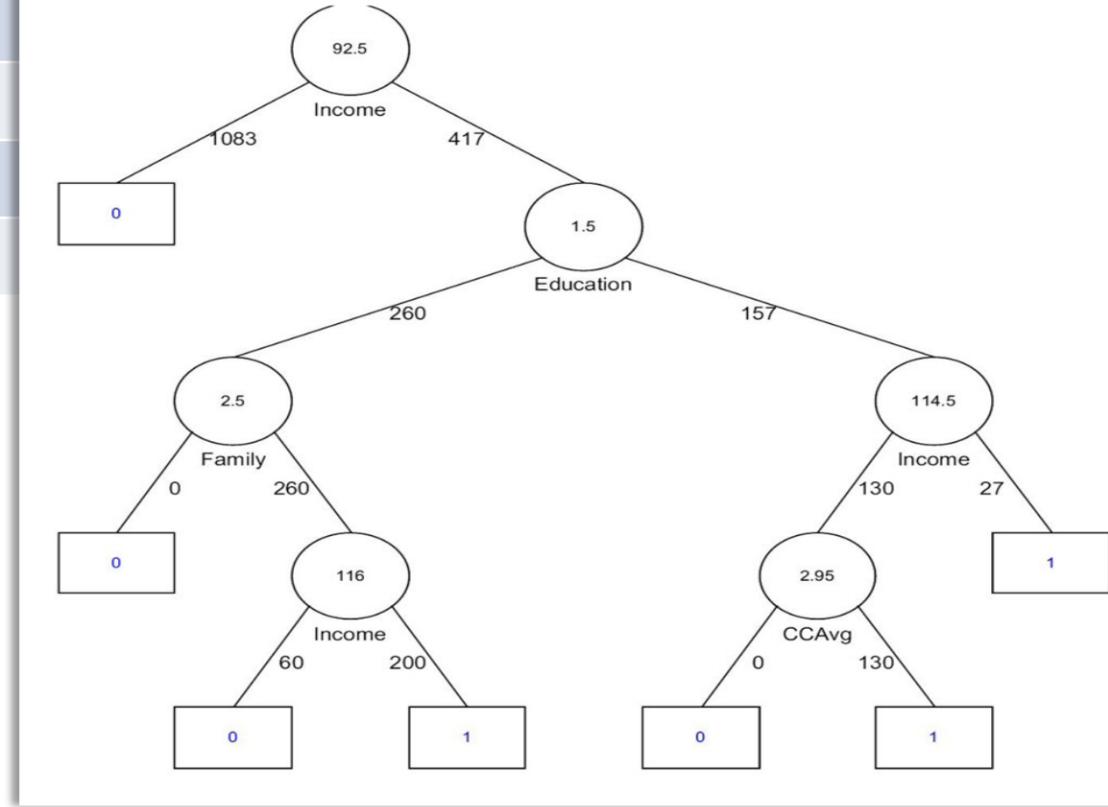
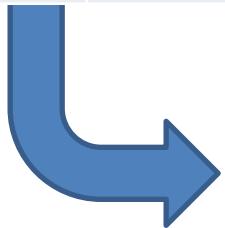
# Tree?





## Training Partition

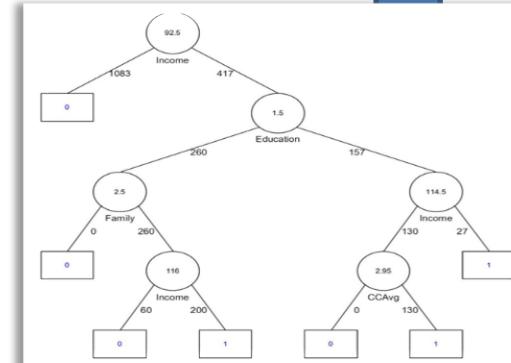
Income	Education	Family Sz	CC Avg	Acceptor?
...	...	...		
...	...	...		
...	...	...		
...	...	...		

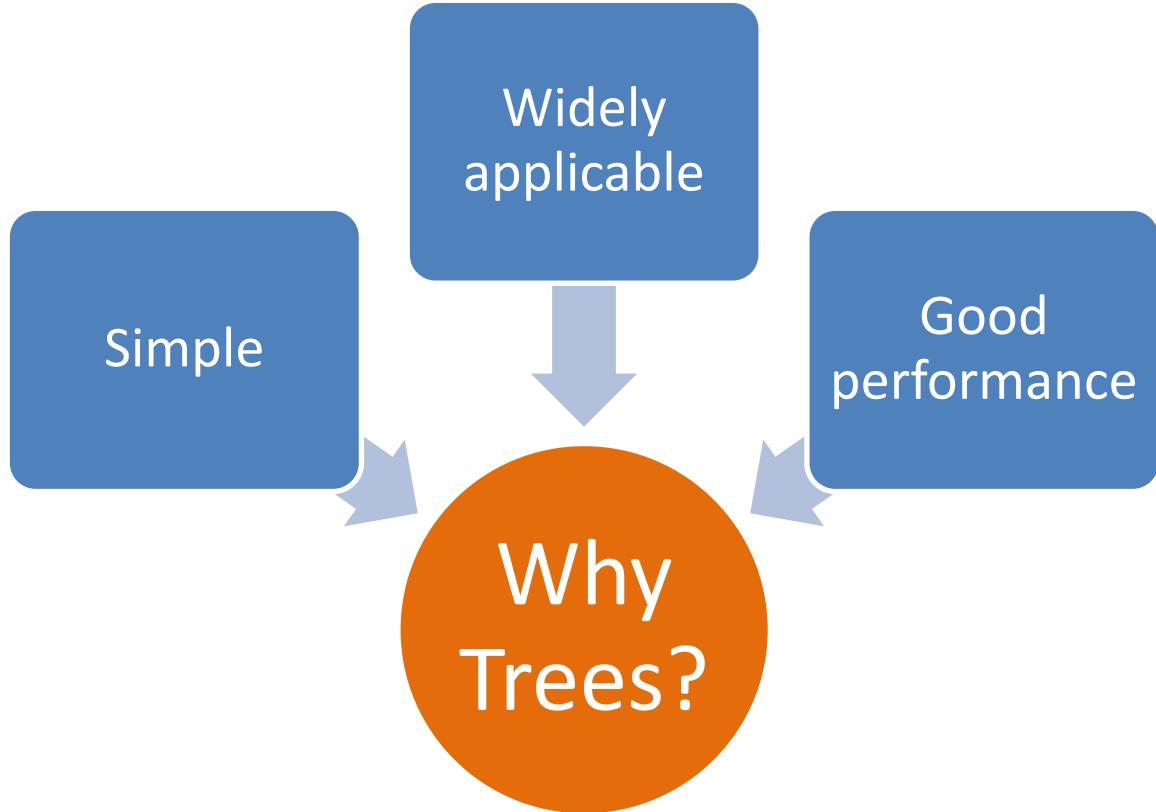


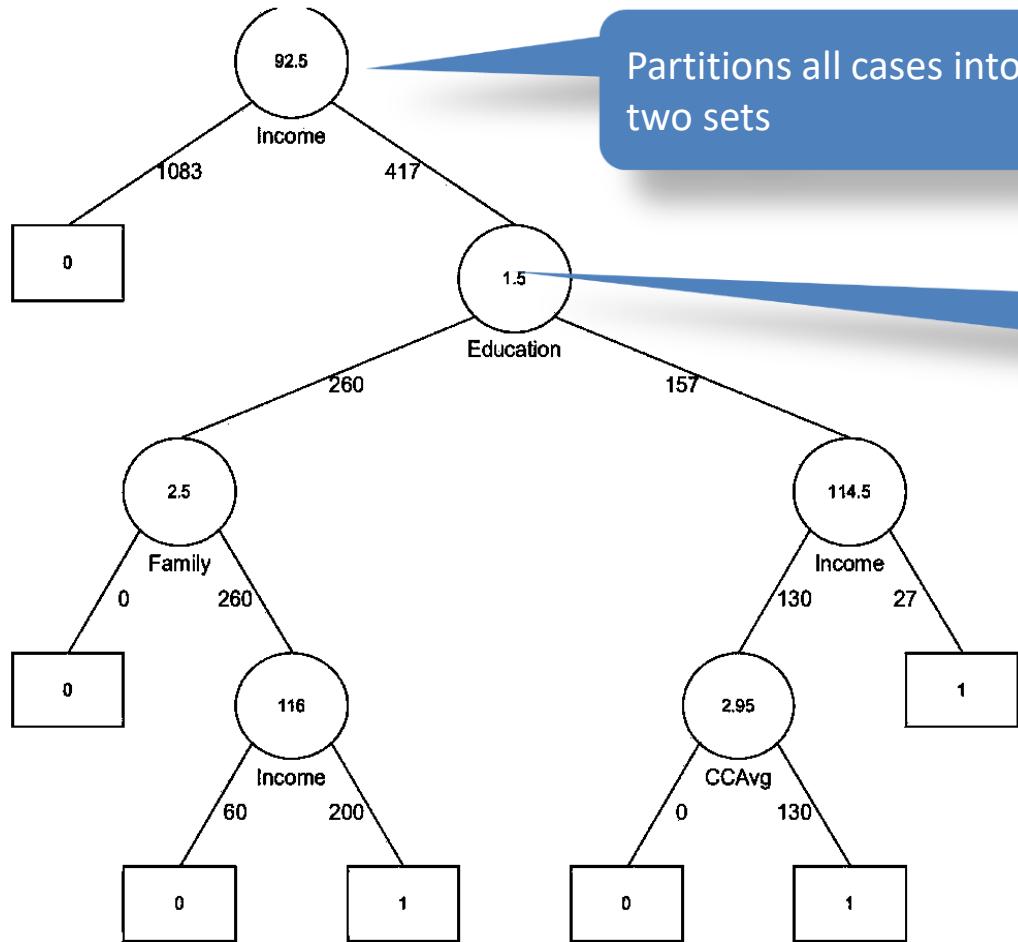
## Test Data

Income	Education	Family Sz	CC Avg	Acceptor?	Model says
...	...	...	...	Yes	No
...	...	...	...	No	Yes
...	...	...	...	No	No
...	...	...	...	...	...

How good is the model?





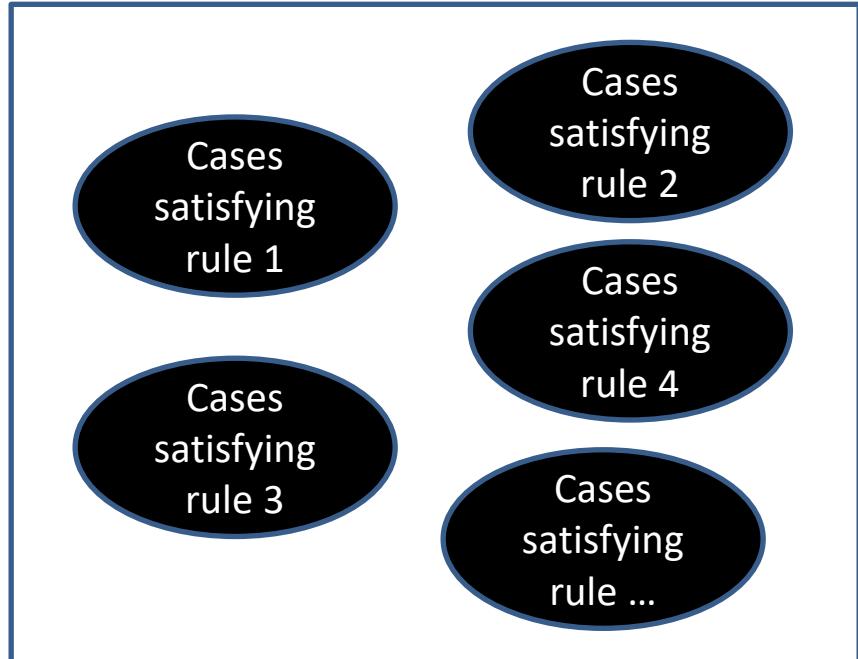
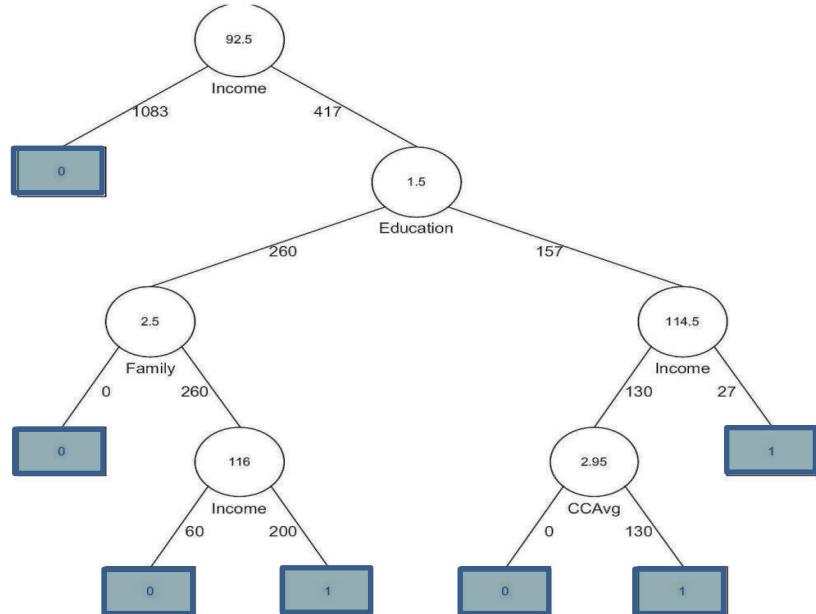


Partitions all cases into  
two sets

Partitions all cases with  
income > 92.5 into two sets

Each “decision node” partitions  
its base cases into two sets based  
on some criterion

# What does a set of rules do?



A complete set of rules **divides** all the cases into disjoint sets such that each set of cases matches a **single** rule.

The above statement is not entirely accurate. It represents an idealized situation that will seldom be applicable in the real world.

It might initially appear that we can always achieve a set of rules that provide perfect matches – after all can we not create a rule for every case? If we do that we will have a huge number of rules, but they will all be perfect. Does this seem to be valid?

In reality the above approach will not generally be valid. We could have two cases with identical values for all the predictor variables, but having different values for the target variable. Concretely, we might have two cases with the same values for income, education, family size and credit card average and yet one of them accepted the loan offer and the other did not. In this case, no matter what rule matches these two cases, the rule can **never** be correct for both cases because they differ in the value of the target variable.

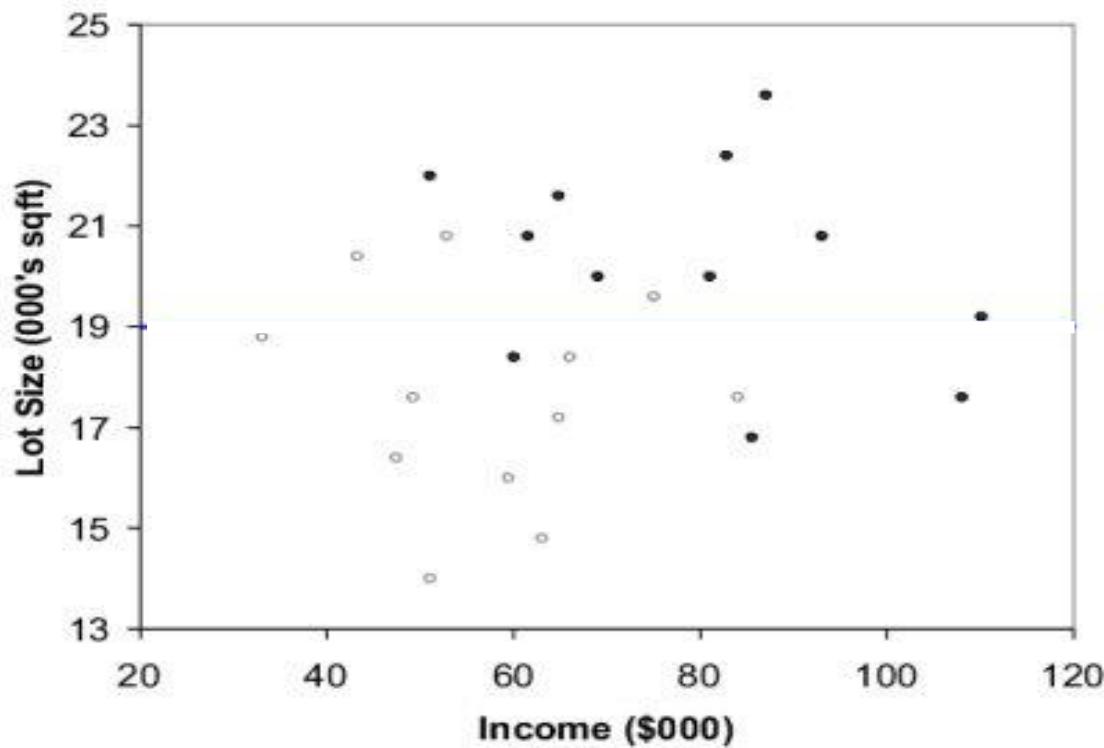
From the above example, we can see clearly that we will not pursue perfect rules. In machine learning we should not expect perfection. Instead we are looking for patterns that apply in enough cases to be of some use.

We are looking for **lift** – that is to beat the averages by a reasonable margin so that applying our models is better than doing things randomly or using more or less obvious ideas.

We see that the set of rules partitions the cases, with each partition consisting of the result of applying a specific rule.

Although we cannot strive for perfection, we still want the partitions to be as homogeneous as possible. Suppose that a particular rule applies to 100 cases in our data set and these cases belong to several different target groups (say, many acceptors and many non-acceptors as well), then the rule might not be as good as another that applies to a majority of a single target group – for example, most of the applicable cases are acceptors. In this case, the rule is perhaps doing a good job of finding an underlying pattern.

# Ownership of riding mower example



This previous picture shows the all the data points. Owners are shown with dark dots.

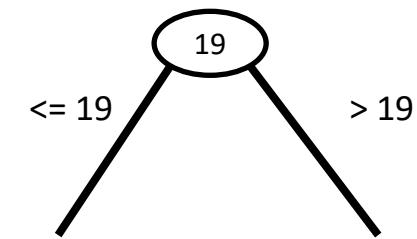
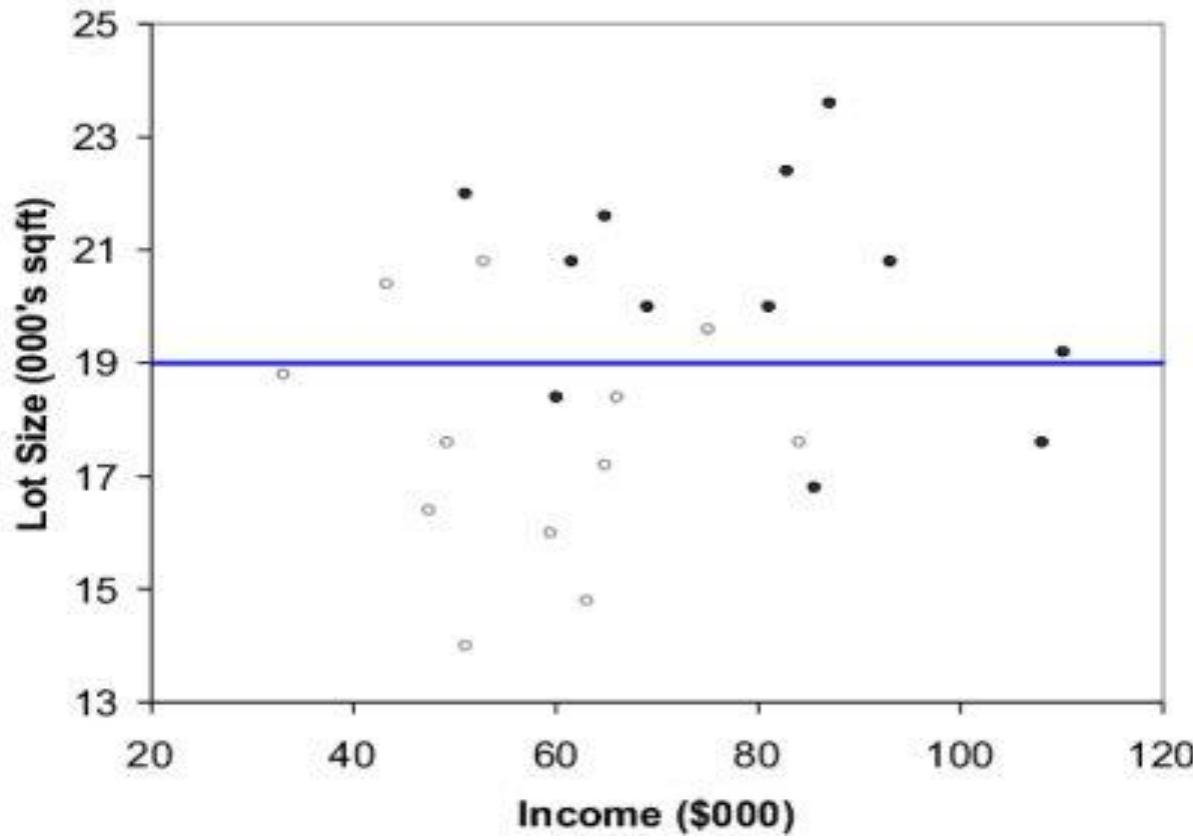
As we observed earlier, while owners generally have high incomes and lot sizes, there are owners with somewhat low incomes and somewhat low lot sizes too.

We have seen that rules effectively partition the data set into a sets of disjoint cases with each set being predominantly homogeneous (say acceptors or non-acceptors). Our procedure for finding rules therefore revolves around building these homogeneous partitions. Once we have that, the rules simply follow.

We want to divide this whole space into regions that contain predominantly owners or non-owners, and derive rules from there.

This procedure is called “**recursive partitioning**”

# Partitioning



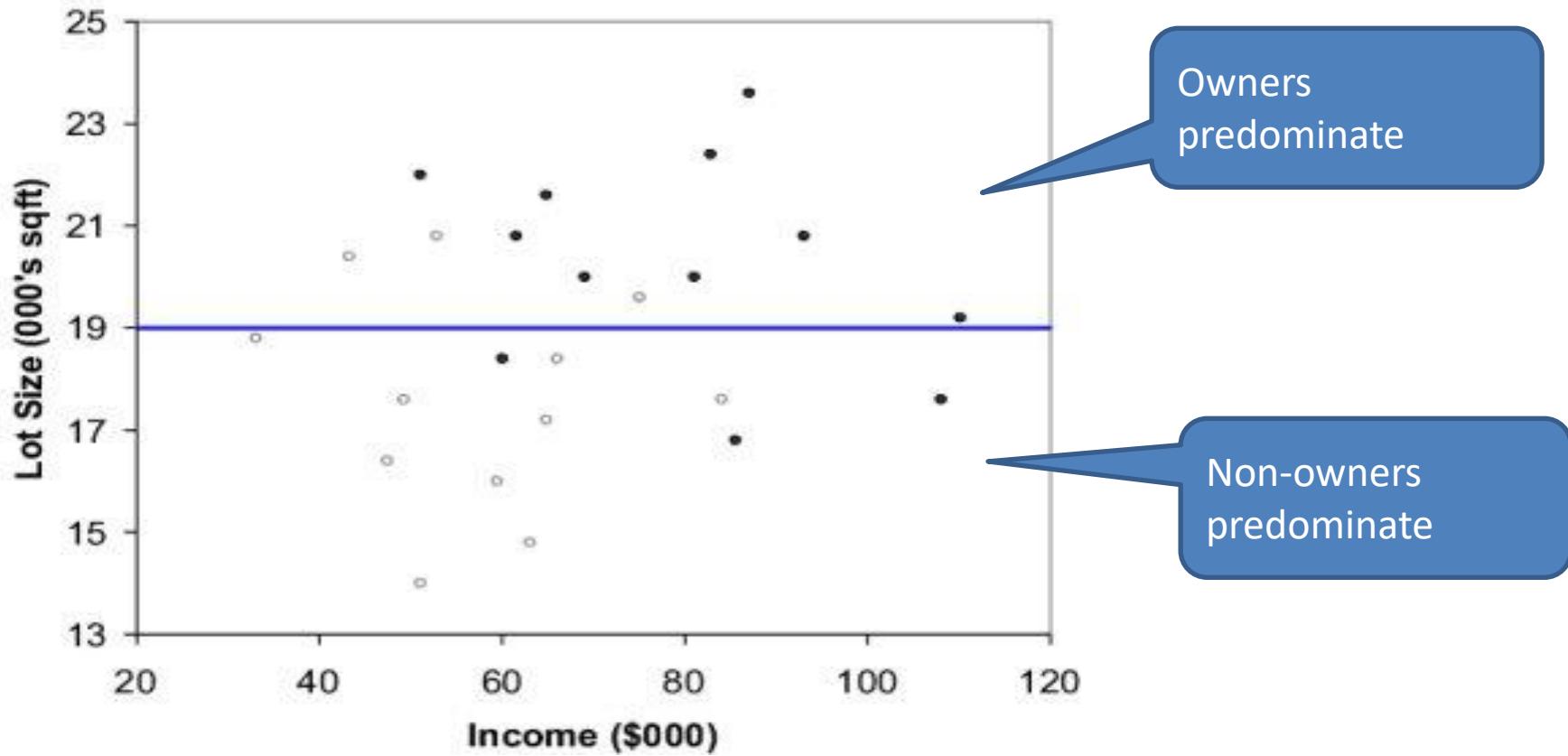
In the above case, our overall region is a rectangle (because we have two variables and therefore have a two-dimensional space) and we have shown above a partition based on lot size. We have shown a line which divides the whole region into two parts. All values above 19 are in one partition and all values  $\leq 19$  are in the other.

We could have chosen to partition by income rather than lot size. Furthermore, having chosen lot size, we could have partitioned at several places.

Why lot size? Why 19?

Before we look at those issues, take a look at the figure – do you think that the partition is “good?” What would be your criterion for deeming a partition to be “good?”

# “Good” Partition Decision?

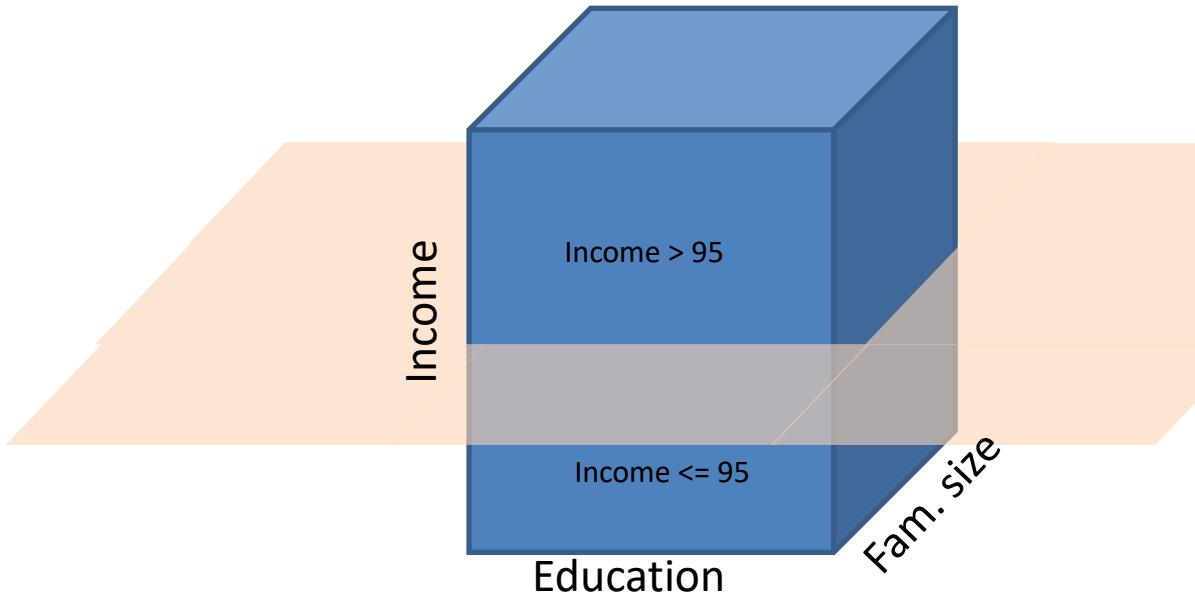


Let us think back to what we are trying to achieve with partitioning – we are trying to create regions with a lot of homogeneity. That is, we want to create regions that have predominantly owners or non-owners.

The partition at 19 seem like a good one because the resulting two regions are much more homogeneous than the original partitioned space was. We now have a majority of owners in the top region and a majority of non-owners in the bottom partition. Visually observe that drawing the partition at any other value of lot size would not have helped us to do better because the two resulting regions would have been less homogeneous.

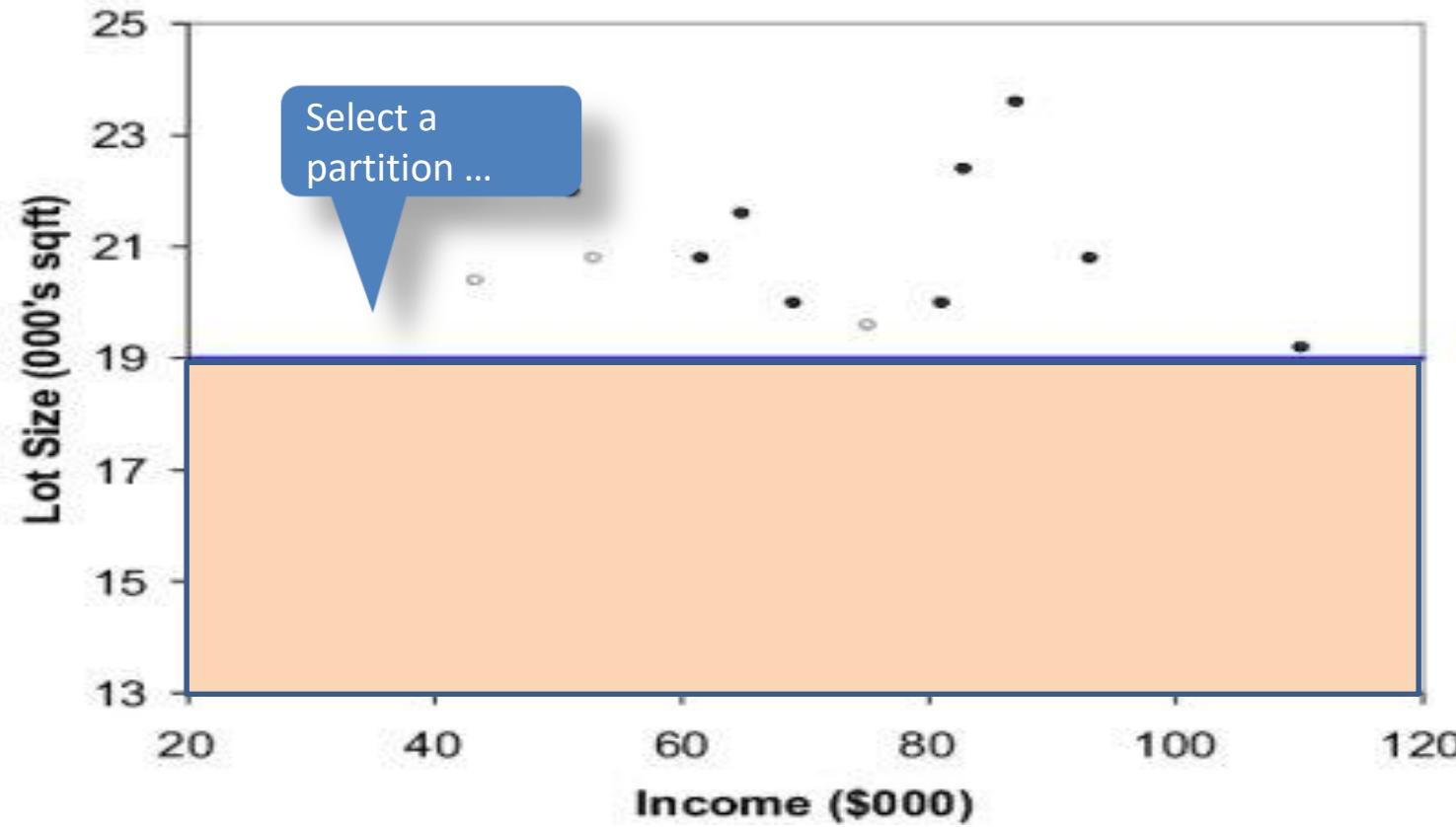
If we had three variables instead of two, then each partition will be a three dimensional space and partitioning is achieved by means of a plane as shown in the next slide.

# Partitioning in Three Dimensions

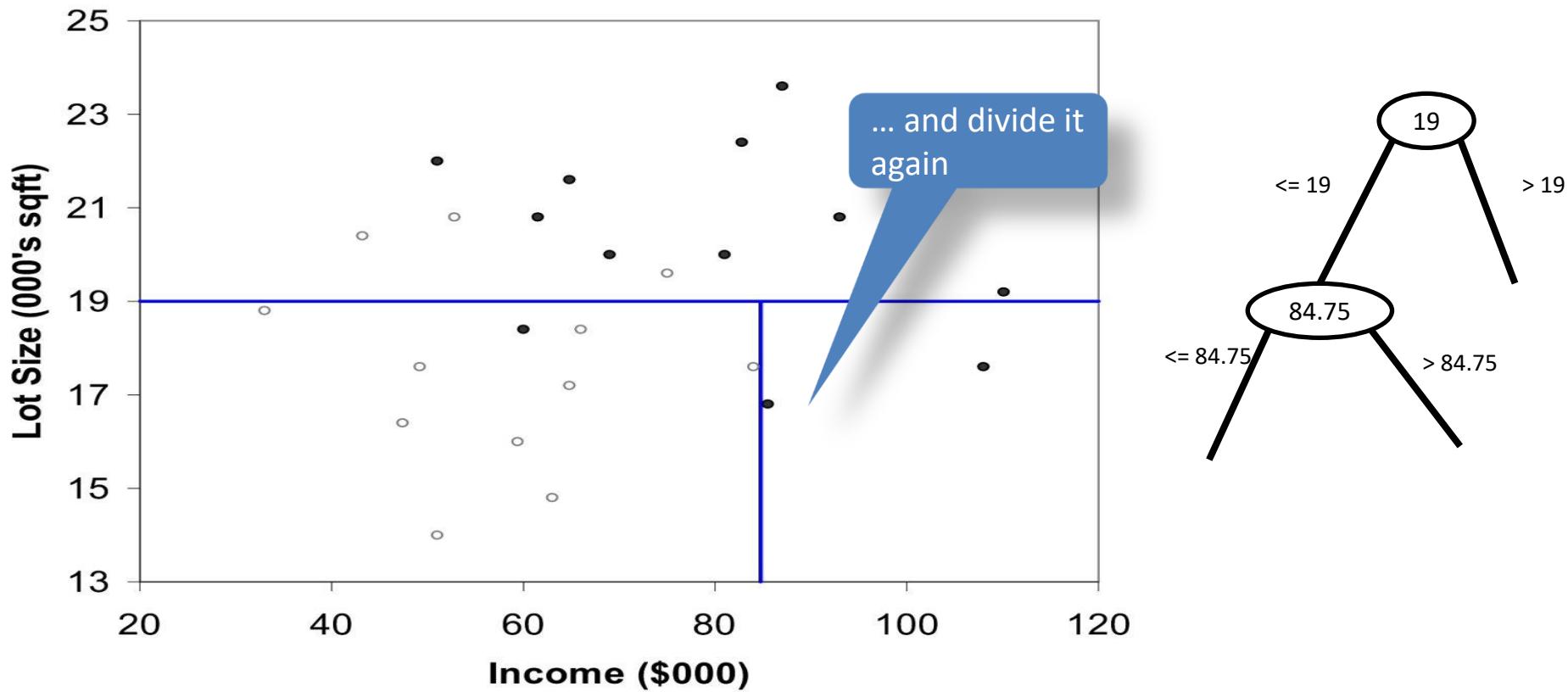


Plane to partition a three dimensional regions into two regions

# Recursive Partitioning



# Recursive Partitioning



... and again divided into two, this time based on income. Those cases with income less than or equal to 92.5 went to one partition and the rest went into another.

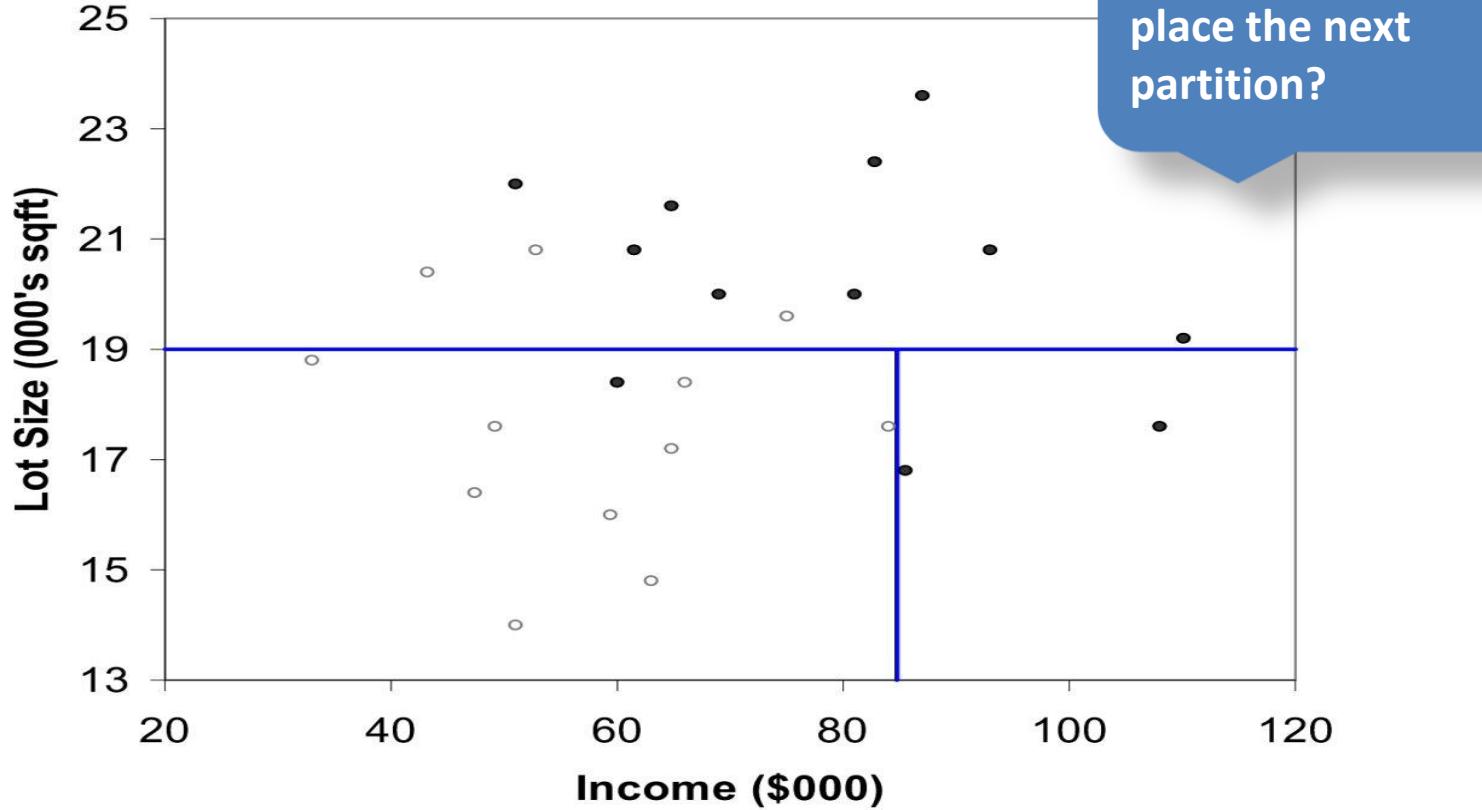
Note however that we are only partitioning the selected region. The top portion remains as is for now.

Several big questions remain:

- Why did I pick the variable lot size selected for the prior partition?
- Why did I choose to partition at 19 and not some other value?
- Why did we select the bottom region?
- Why did we select Income as the variable for partitioning?
- Why did we select 92.5 as the specific value to partition on?

We have already hinted that “homogeneity” of the regions resulting from a partitioning step guides our choices.

# Your Turn



Where would you place the next partition?

Take a look at the situation and consider where you might choose to place the next partition. Remember, you have to select a region (from among the three above) and also select where exactly to draw the line.

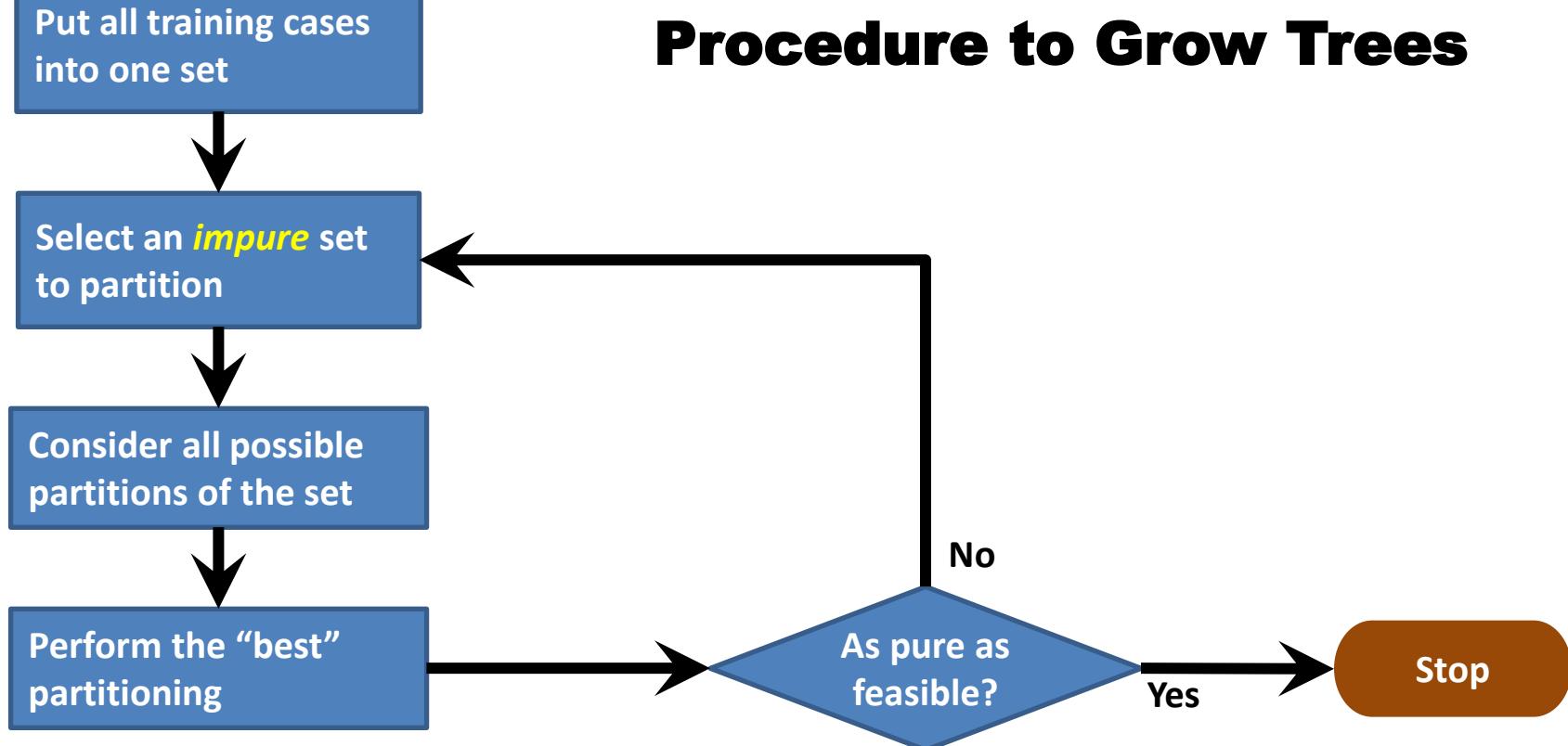
For example, suppose you chose the topmost region, you then have to decide whether to partition on lot size or income. If you choose lot size, then you will select the value to partition on. Again suppose you choose 23 (just for illustration only, might not be a good choice), then you are effectively drawing a horizontal line running through the selected region at 23 on the Y axis

Perfection not always possible.

Under what conditions might obtaining perfect partitioning be impossible?

Answer: If two cases have the same values for the predictor variables and differ only in their values for the target variable, then getting a set of rules that performs perfectly is impossible.

# Procedure to Grow Trees



We want each group to consist only of one type of case (Owner or Non-owner).

Therefore, we need a measure of the extent to which a group is impure (or lacks homogeneity). Clearly, if all the cases in a group are owners or all are non-owners, then the group is pure. On the other hand the maximum amount of impurity occurs when the group is perfectly heterogeneous. In the present example that would mean that the group consists of 50% owners and 50% non-owners. This is the maximum possible deviation that we can get from absolute purity when every member is of the same type. (A 60-40 or 70-30 mix of owners and non-owners is purer than a 50-50 mix because there is a tendency for the predominance of some group.)

When we partition, we change the measure of impurity.

At each stage, among all possible partitions, we select the one that results in the best improvement in the measure of purity. This is the same as saying that we choose the split that results in the best possible improvement in overall purity.

# Homogeneity

- The main goal in a decision tree algorithm is to identify a variable and classification that results in a more homogeneous distribution with reference to the target variable
- The homogeneous distribution means that similar values of the target variable are grouped together so that a concrete decision can be made

# Entropy

A measure of the randomness in the information being processed. The higher the entropy, the harder it is to draw any conclusions from that information.

$$H(X) = - \sum_{i=1}^n P(x_i) \log_b P(x_i)$$

Where H is the entropy. X is the dataset. P(x) is the proportion of x in X, n is the number of classes in the target variable

# Information Gain

The amount of information gained about a random variable from another random variable. The difference between the entropy of *parent* node and weighted average entropy of *child* nodes.

$$IG(S, A) = H(S) - H(S, A)$$

Alternatively,

$$IG(S, A) = H(S) - \sum_{i=0}^n P(x) * H(x)$$

Where  $IG(S, A)$  is the difference in entropy.

$H(S)$  is the entropy of  $S$

$H(S, A)$  is the entropy of  $S$  given  $A$ .

# Gini Index: A Measure of Impurity of a Set

$m$ : number of classes



2 – Owner and Non-Owner

$p_k$ : proportion of elements in class  $k$



12 Owners and 12 Non-Owners  
 $p_1$  and  $p_2$  equal 0.5

$$Gini(x) = 1 - \sum_{k=1}^m p_k^2$$

$$0 \text{ to } \frac{(m-1)}{m}$$

# Gini Calculation for Initial Set

m

2 – Owner and Non-Owner

$p_1$

0.5: 12 Owners out of 24

$p_2$

0.5: 12 Non-owners out of 24

$$Gini(x) = 1 - \sum_{k=1}^m p_k^2$$

$$1 - (0.5^2 + 0.5^2) = 0.5$$

# Entropy: Another Impurity Measure of a set

m: number of classes

2 – Owner and Non-Owner

$p_k$ : proportion of elements in class k

12 Owners and 12 Non-Owners  
 $p_1$  and  $p_2$  equal 0.5

$$E(s) = - \sum_{k=1}^m p_k \log_2 p_k$$

0 to 1

# Entropy Calculation for Initial Set

m

2 – Owner and Non-Owner

$p_1$

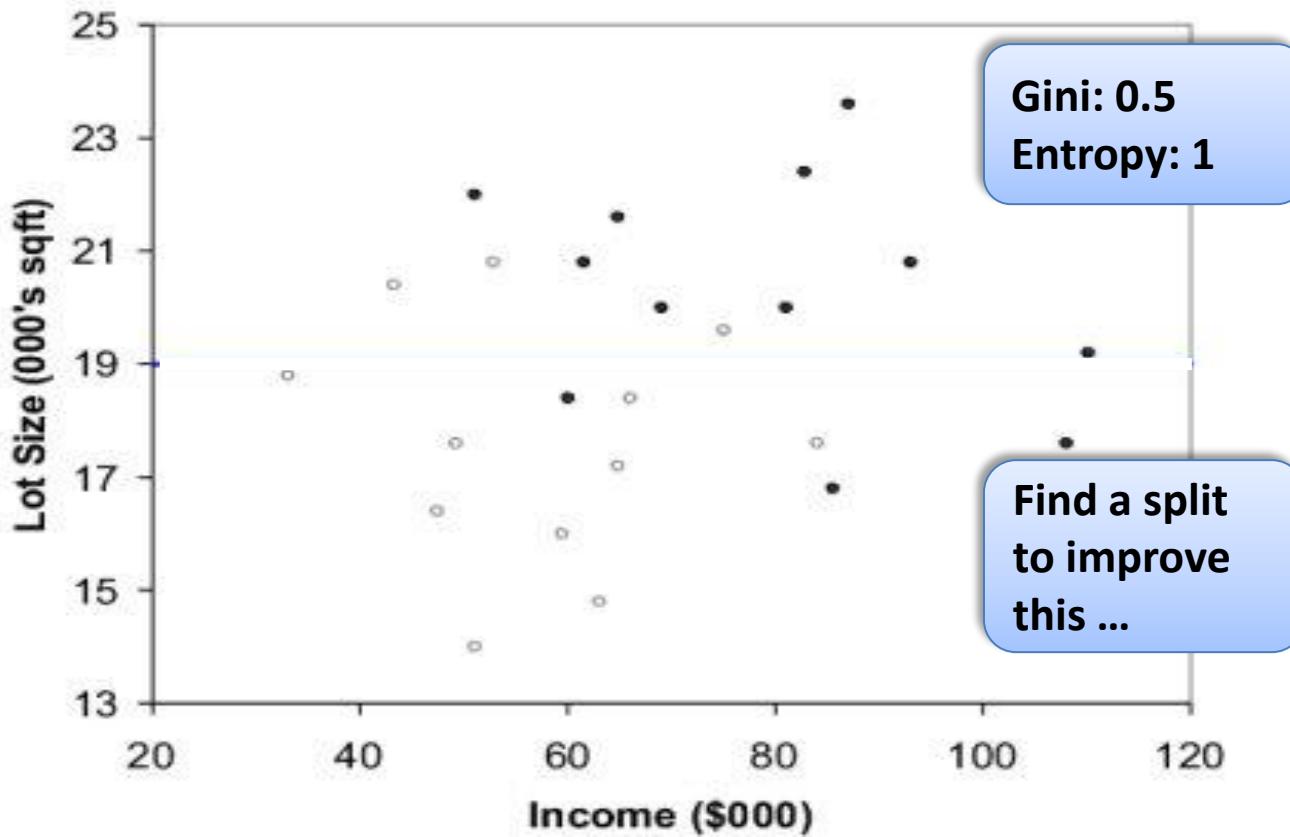
0.5: 12 Owners out of 24

$p_2$

0.5: 12 Non-owners out of 24

$$E(s) = - \sum_{k=1}^m p_k \log_2 p_k$$

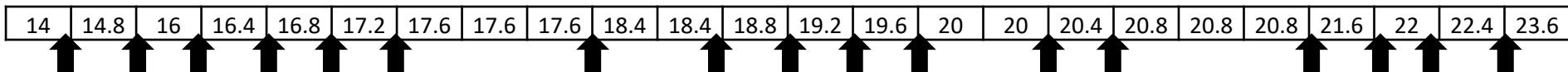
$$-[0.5*(-1)+0.5*(-1)] = 1$$



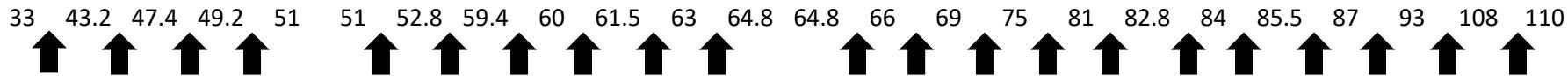
# Finding the First Split

Lot Sizes -- sorted

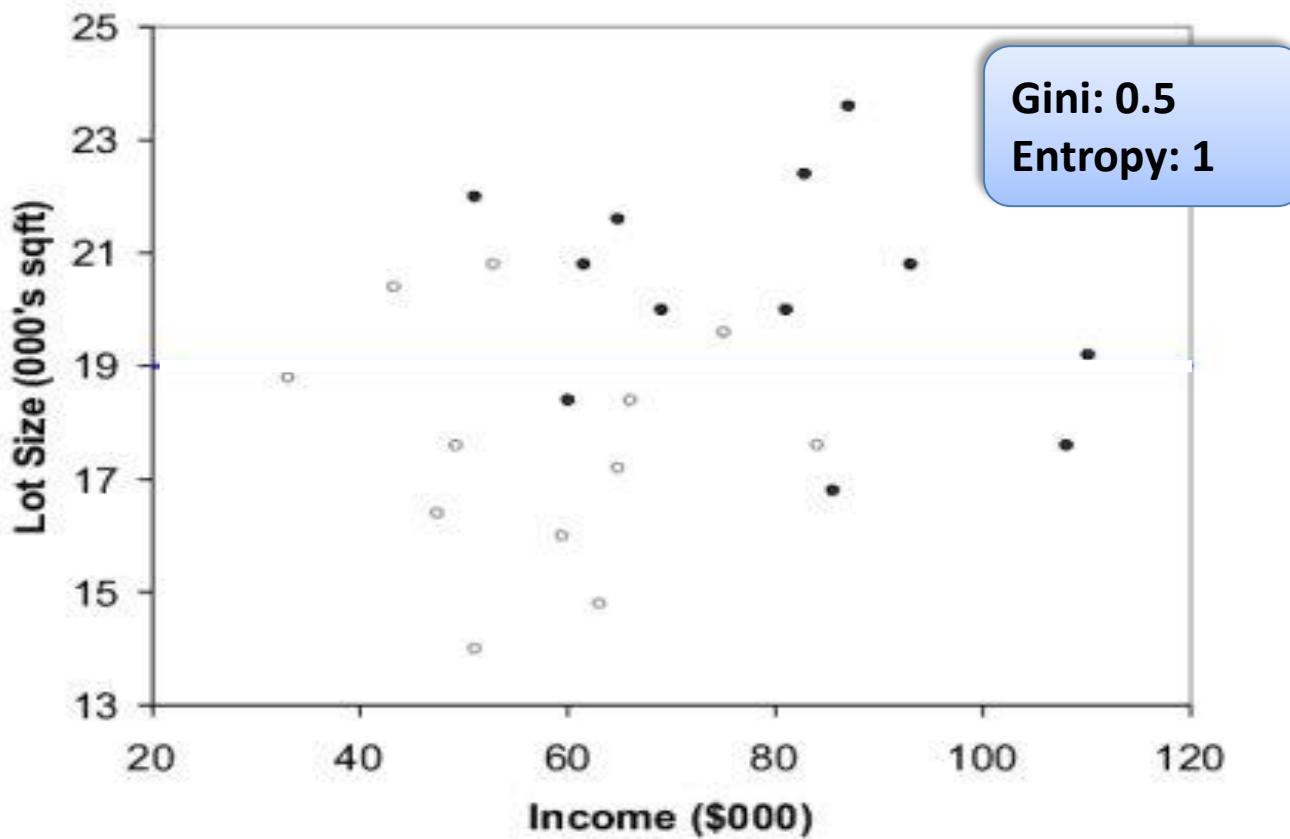
Can split between any two successive values that are different



Incomes -- sorted

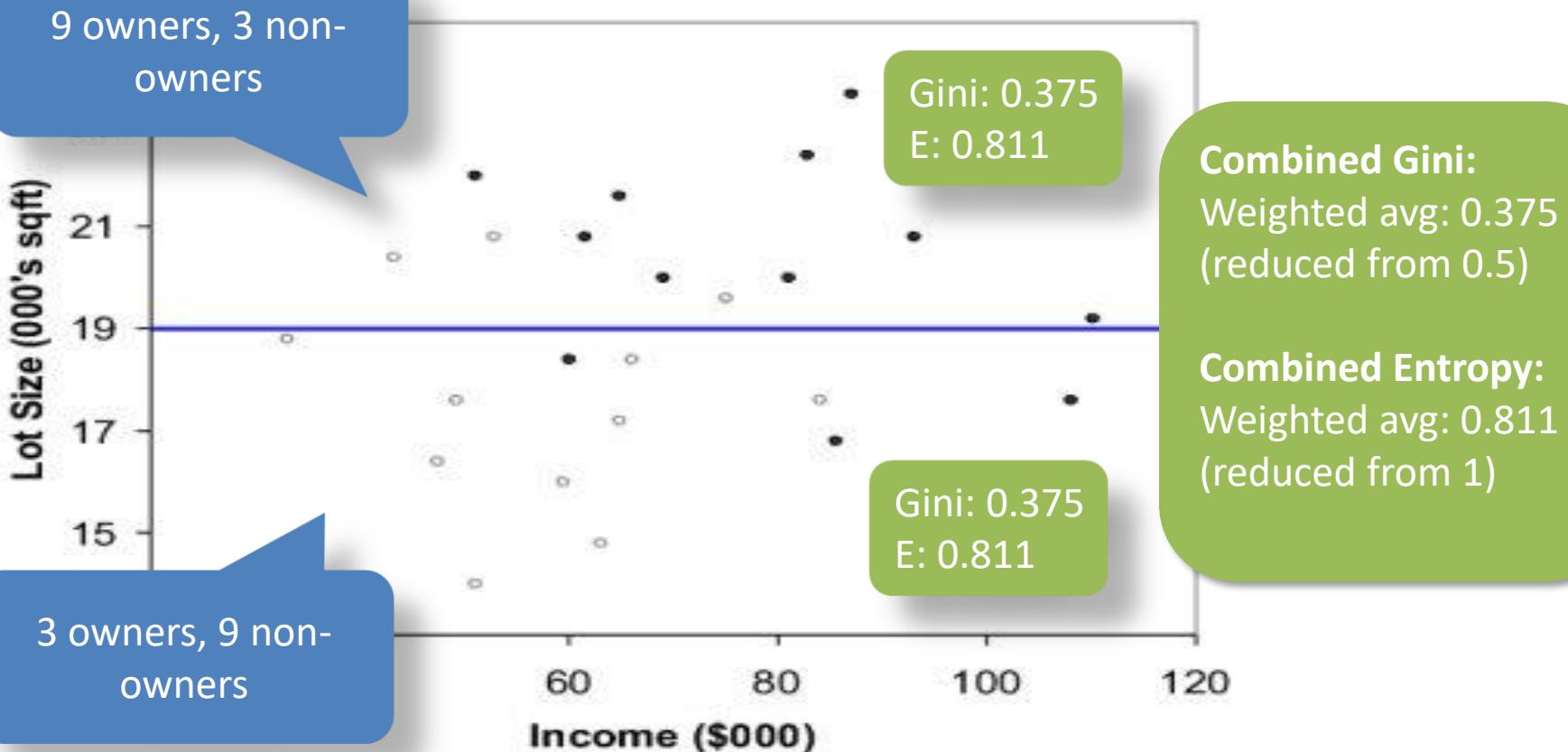


38 -- Try all and pick the one that reduces overall impurity the most

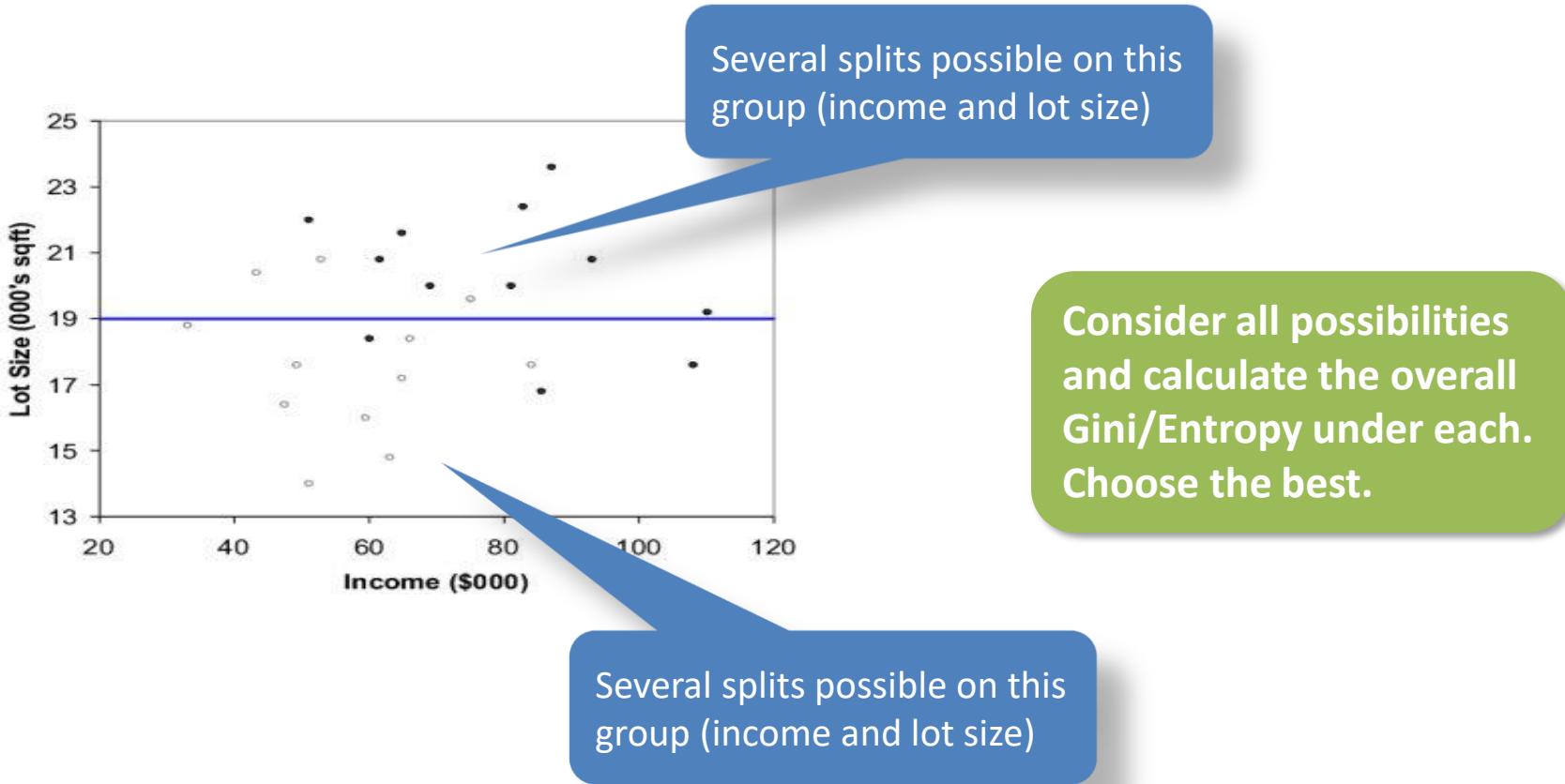


Gini: 0.5  
Entropy: 1

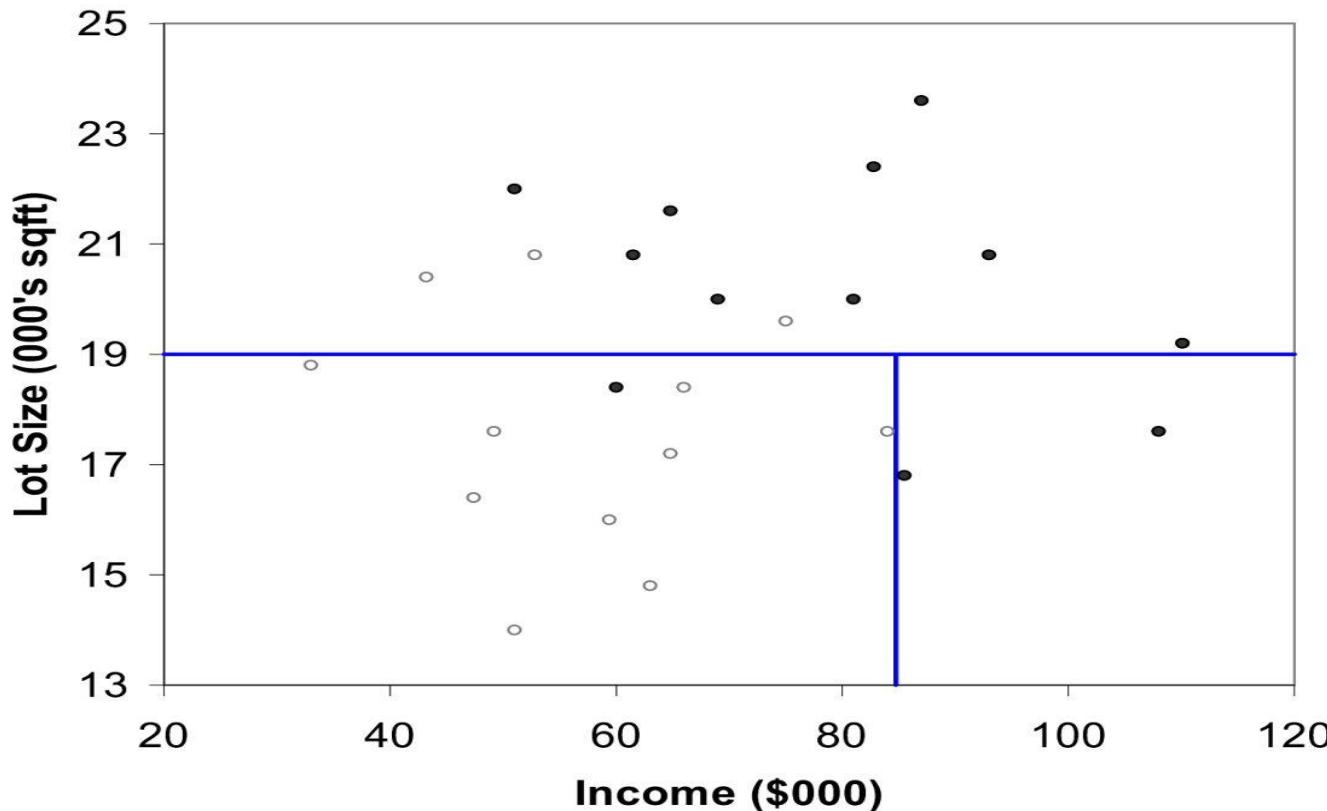
# First Split – Lot Size 19



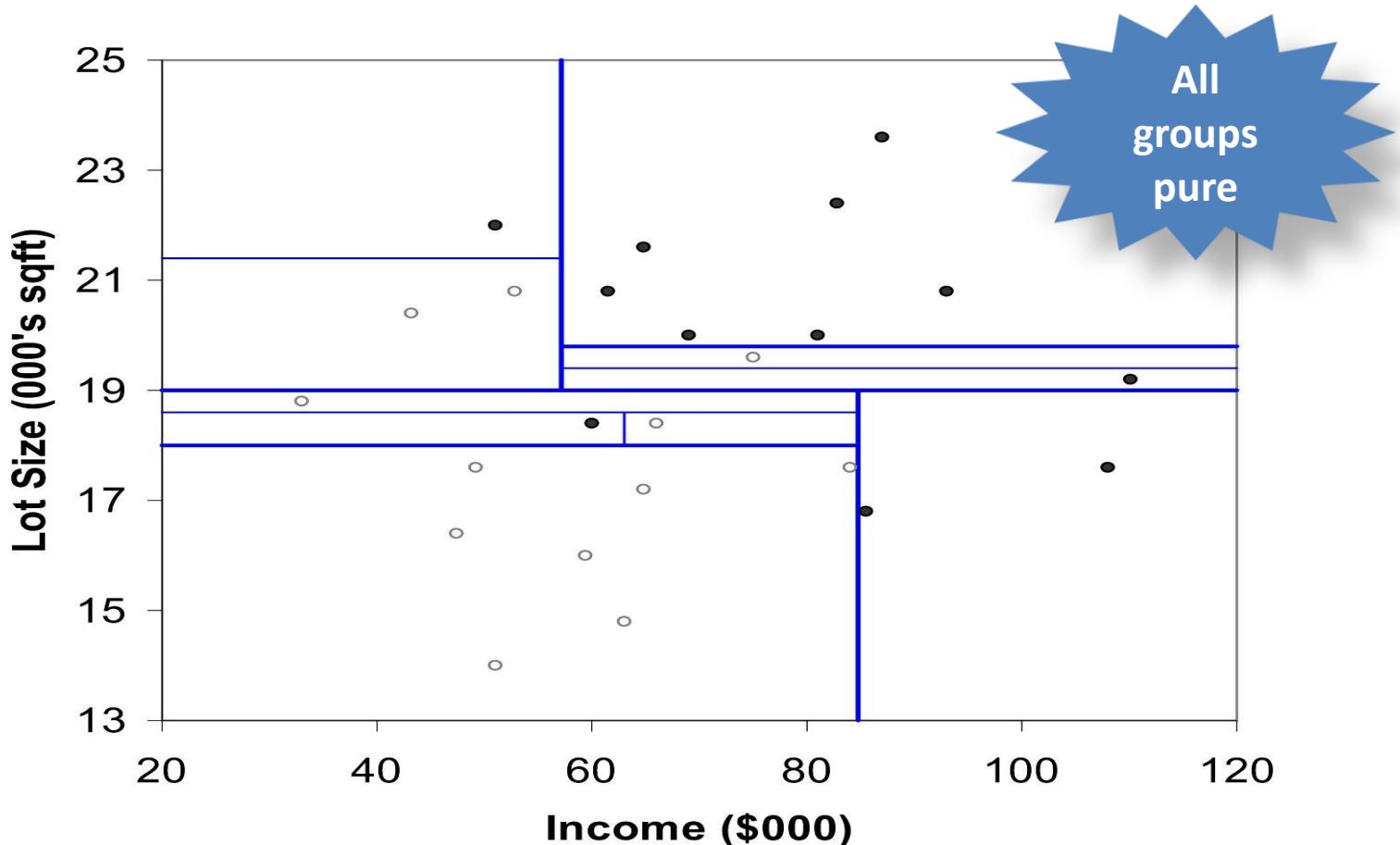
# Choosing Subsequent Splits

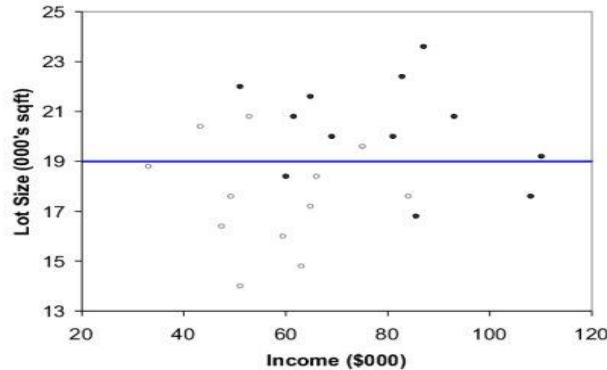


# Second Split – Income \$84K

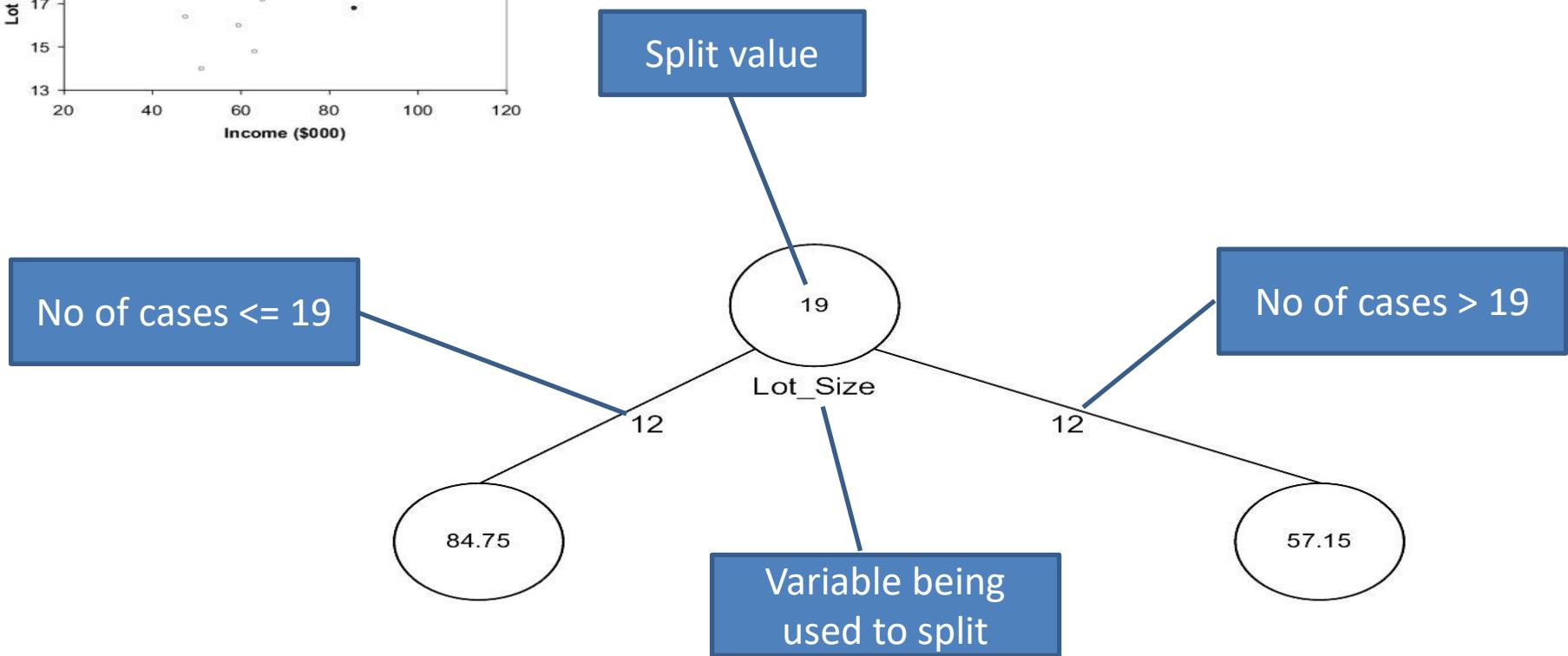


# After All Splits

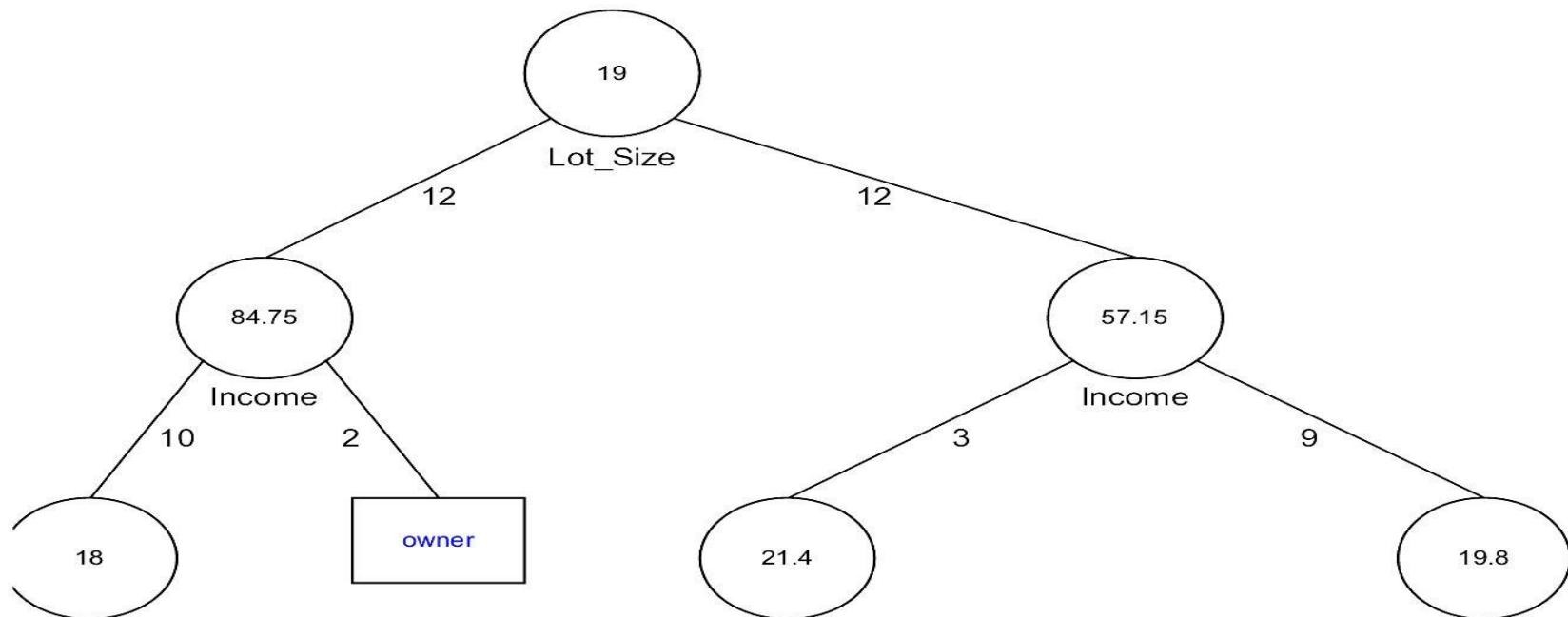


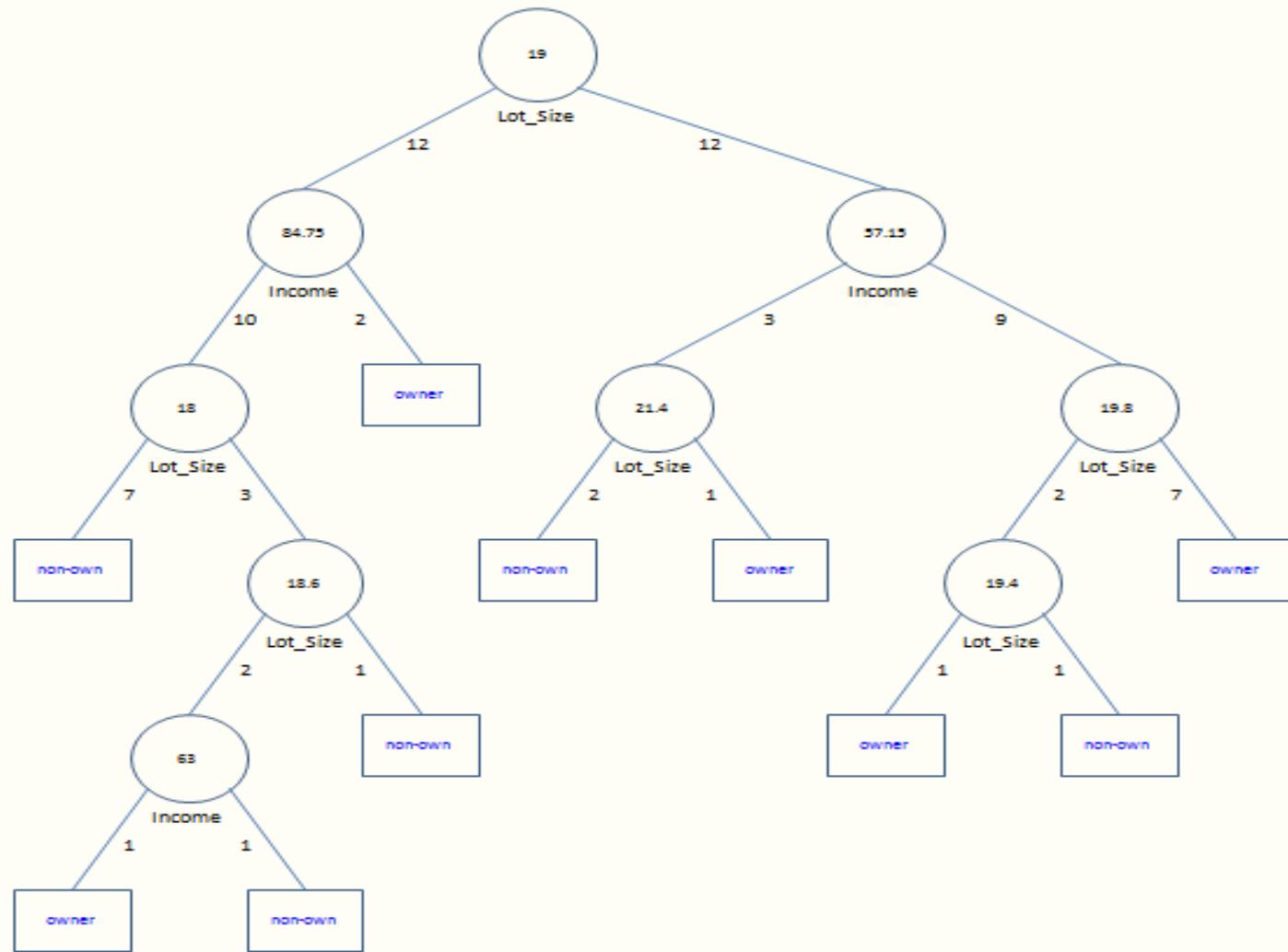


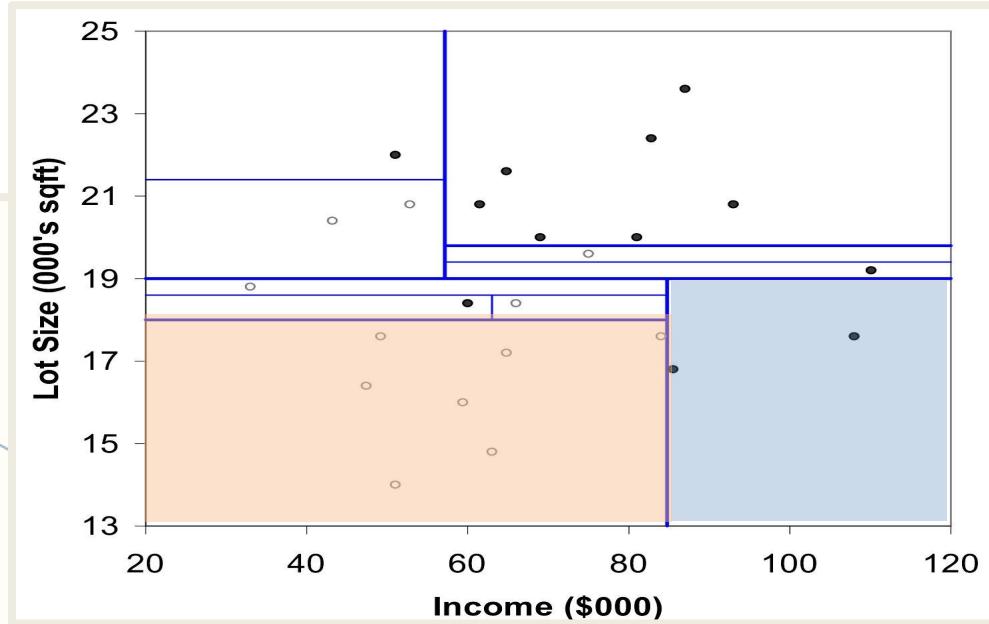
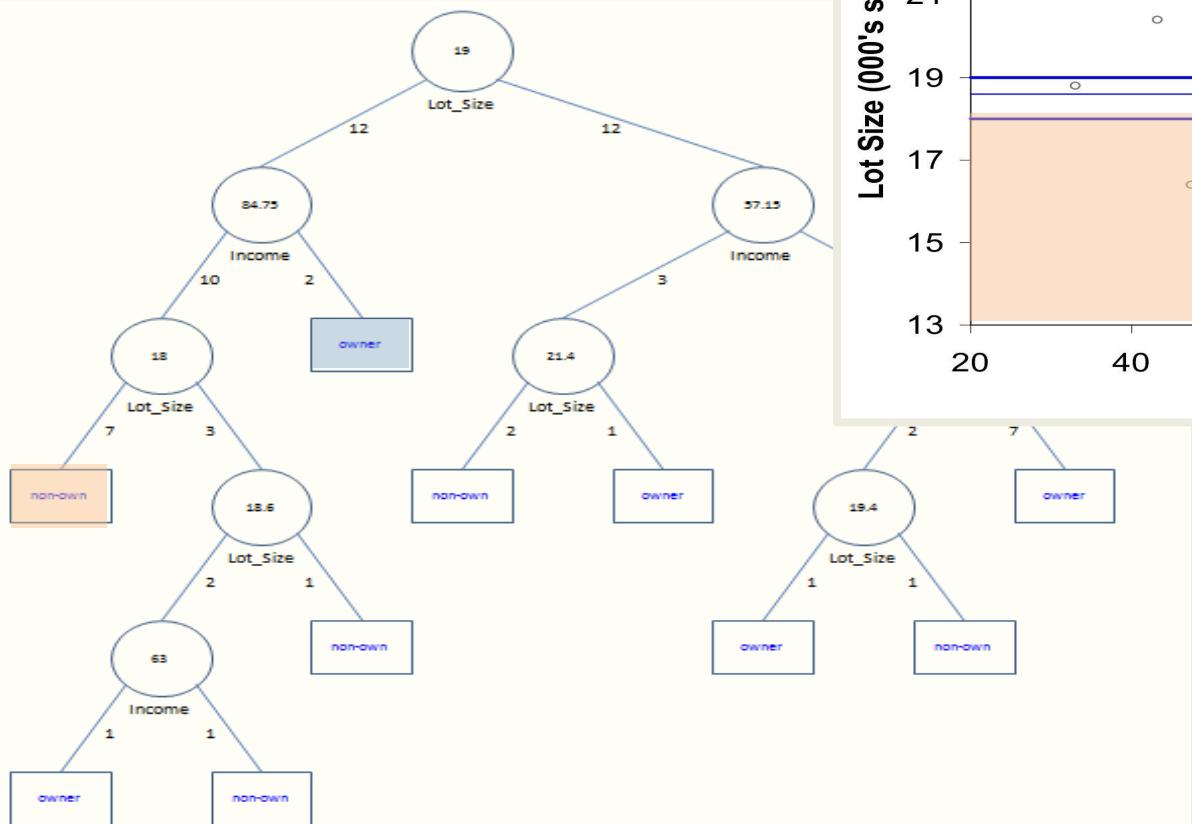
# Tree Representation



# Tree -- After 3 Splits

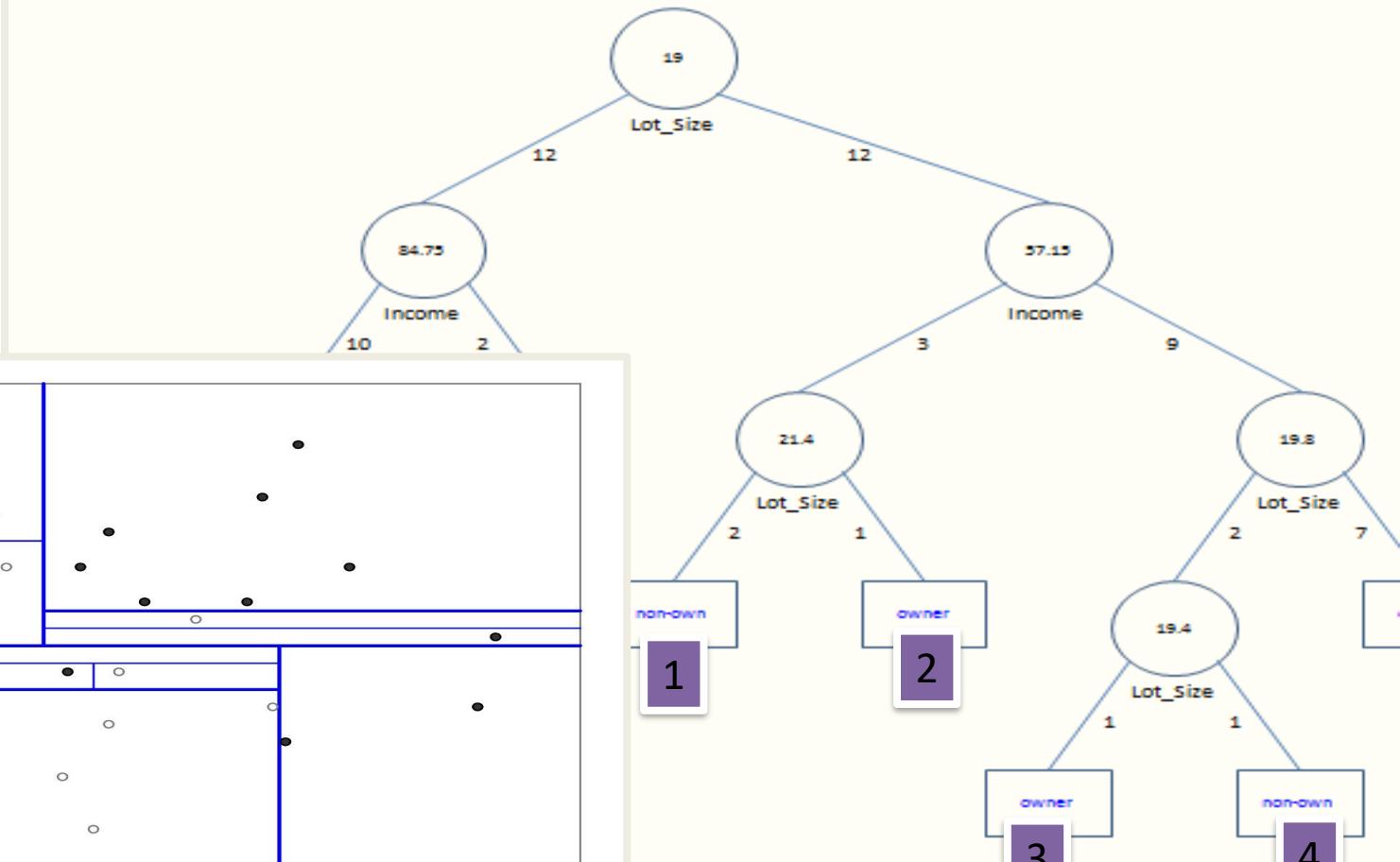
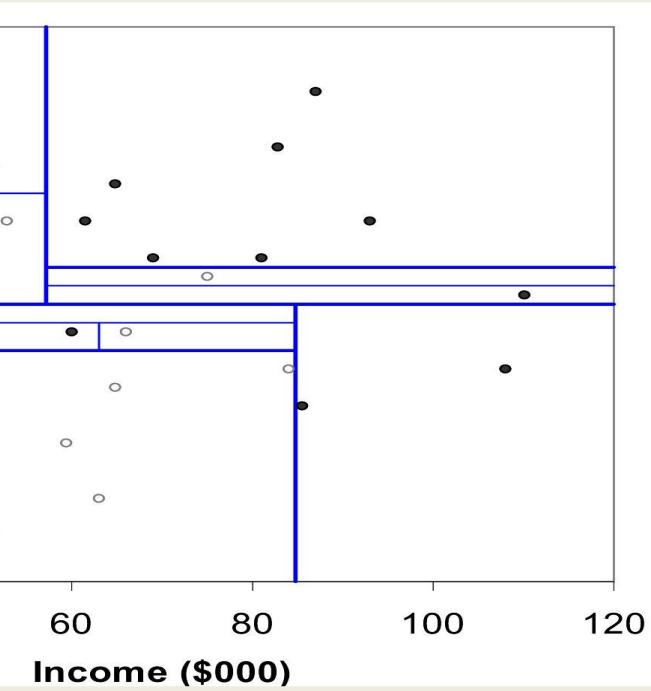




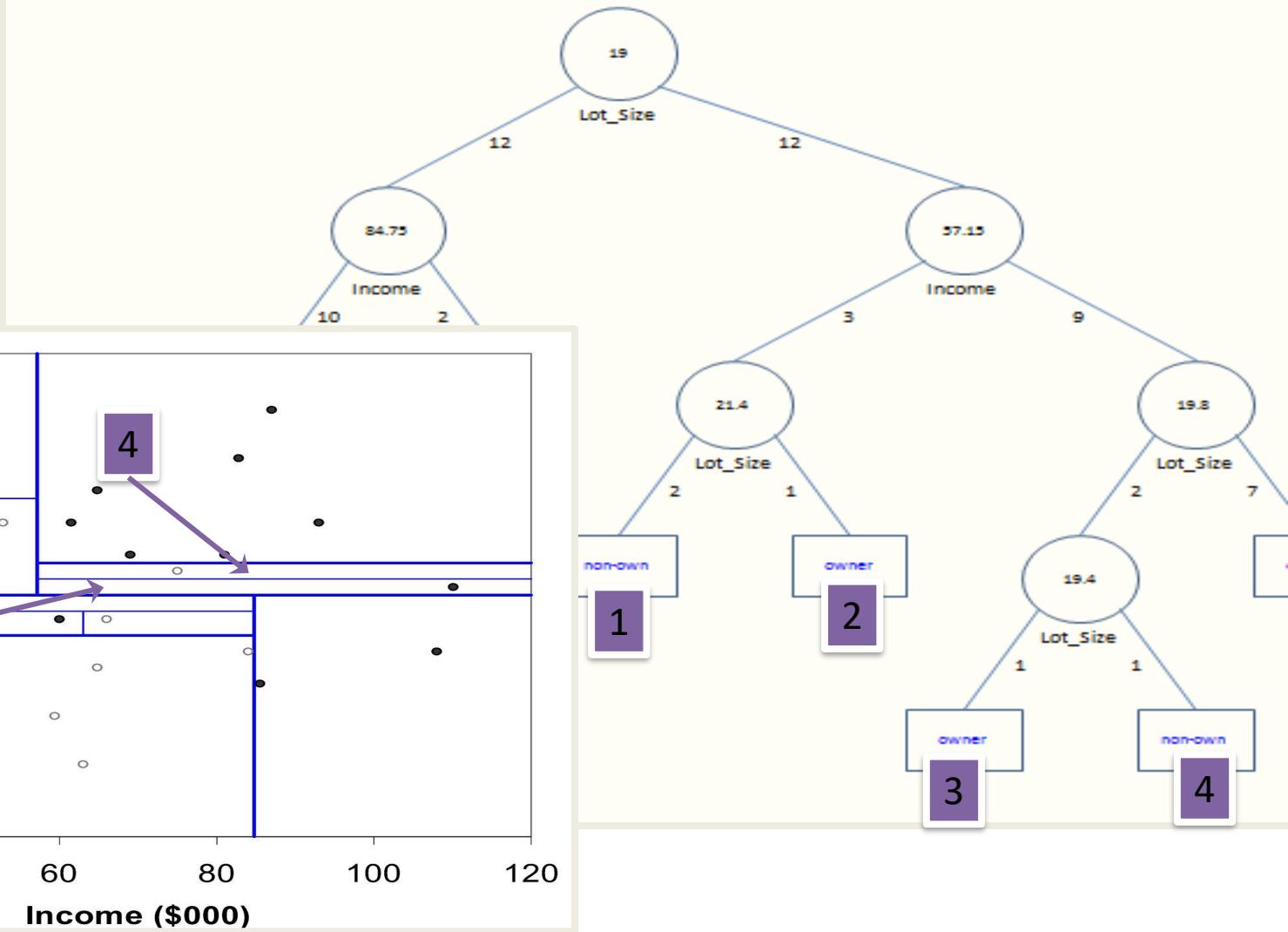


# After All Splits

# Your Turn



# Answer



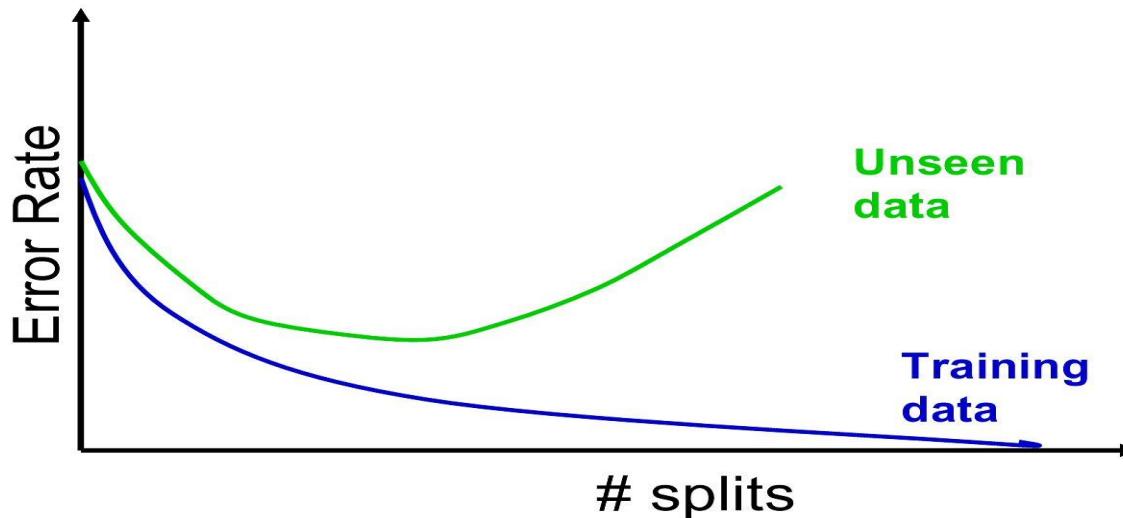


100%  
purity

Over-  
fitting

Beyond a point, fitting noise -- Over fitting  
Reduces performance on validation data

# Impact of Number of Splits



# Iterative Dichotomizer 3 (ID3) algorithm steps (uses entropy)

1. Calculate the initial *entropy* of the system based on the target variable.
2. Calculate the *information gains* for each candidate variable for a node. Select the variable that provides the maximum information gain as a decision node.
3. Repeat step 2 for each branch (value) of the node (variable) identified in step 2. The newly identified node is termed as *leaf* node.
4. Check whether the leaf node classifies the entire data perfectly. If not, repeat the steps from step 2 onwards. If yes, stop.

# ID3 algorithm example

Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes
Rain	Mild	High	Strong	No

# Candidate Variable: *Outlook*

		play		total
		yes	no	
Outlook	sunny	3	2	5
	overcast	4	0	4
	rainy	2	3	5
				14

$$E(S) = -[(9/14)\log(9/14) + (5/14)\log(5/14)] = 0.94$$

$$\begin{aligned} E(S, \text{outlook}) &= (5/14)*E(3,2) + (4/14)*E(4,0) + (5/14)*E(2,3) = \\ &(5/14)(-(3/5)\log(3/5)-(2/5)\log(2/5)) + (4/14)(0) + (5/14) \\ &((2/5)\log(2/5)-(3/5)\log(3/5)) = 0.693 \end{aligned}$$

# Compare Info Gain for all variables

$$IG(S, \text{outlook}) = 0.94 - 0.693 = 0.247$$

$$IG(S, \text{Temperature}) = 0.940 - 0.911 = 0.029$$

$$IG(S, \text{Humidity}) = 0.940 - 0.788 = 0.152$$

$$IG(S, \text{Windy}) = 0.940 - 0.8932 = 0.048$$

Highest Information Gain is for *outlook*. *Outlook* is first decision node

Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes

Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Overcast	Hot	High	Weak	Yes
Overcast	Cool	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes

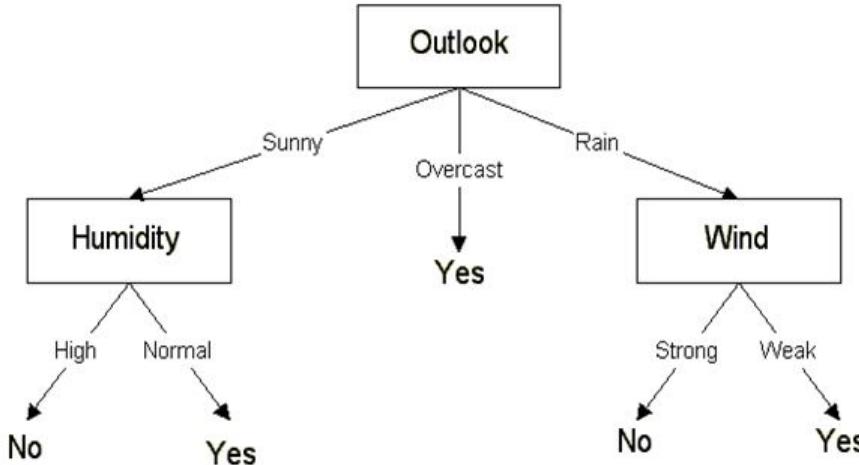
Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Rain	Mild	Normal	Weak	Yes
Rain	Mild	High	Strong	No

Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes

		play			
		yes	no	total	
		hot	0	2	2
Temperature	cool		1	1	2
	mild		1	0	1
					5

$$E(\text{sunny}) = \left(-\frac{3}{5}\log\left(\frac{3}{5}\right) - \frac{2}{5}\log\left(\frac{2}{5}\right)\right) = 0.971.$$

$$E(\text{sunny}, \text{Temperature}) = \frac{2}{5} \cdot E(0,2) + \frac{2}{5} \cdot E(1,1) + \frac{1}{5} \cdot E(1,0) = \frac{2}{5} = 0.4$$



$$IG(\text{sunny}, \text{Temperature}) = 0.971 - 0.4 = 0.571$$

$$IG(\text{sunny}, \text{Humidity}) = 0.971$$

$$IG(\text{sunny}, \text{Windy}) = 0.020$$

Humidity has the highest IG. Continue building out the decision tree

# Gini Impurity

A measure of how often a randomly chosen element from the set would be incorrectly labeled if it was randomly labeled according to the distribution of labels in the subset.

$$Gini(E) = 1 - \sum_{j=1}^c p_j^2$$

# CART algorithm steps (uses Gini)

1. Calculate the *Gini* index based on the target variable.
2. Calculate the *Gini* gains for each candidate variable for a node.  
Select the variable that provides the maximum Gini gain as a decision node.
3. Repeat step 2 for each branch (value) of the node (variable) identified in step 2. The newly identified node is termed as leaf node.
4. Check whether the leaf node classifies the entire data perfectly.  
If not, repeat the steps from step 2 onwards. If yes, stop.

# CART algorithm example

Outlook	Temperature	Humidity	Wind	Played football(yes/no)
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Overcast	Hot	High	Weak	Yes
Rain	Mild	High	Weak	Yes
Rain	Cool	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Overcast	Cool	Normal	Strong	Yes
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes
Rain	Mild	High	Strong	No

		play			
		yes	no	total	
Outlook	sunny	3	2	5	
	overcast	4	0	4	
	rainy	2	3	5	
				14	

$$\text{Gini}(S) = 1 - [(9/14)^2 + (5/14)^2] = 0.4591$$

$$\begin{aligned}
 \text{Gini}(S, \text{outlook}) &= (5/14)\text{gini}(3,2) + (4/14)*\text{gini}(4,0) + (5/14)*\text{gini}(2,3) \\
 &= (5/14)(1 - (3/5)^2 - (2/5)^2) + (4/14)*0 + (5/14)(1 - (2/5)^2 - (3/5)^2) = \\
 &0.171 + 0 + 0.171 = 0.342
 \end{aligned}$$

Gini gain (S, outlook) = 0.459 - 0.342 = 0.117

Gini gain(S, Temperature) = 0.459 - 0.4405 = 0.0185

Gini gain(S, Humidity) = 0.459 - 0.3674 = 0.0916

Gini gain(S, windy) = 0.459 - 0.4286 = 0.0304

Highest Gini Gain is for *outlook*. *Outlook* is first decision node

# **Handling numerical variables**

1. Sort the dataset based on the numerical variable in ascending order.
2. Mark the numerical variable ranges where the target variable transitions from one category to another.
3. Create a separate categorical variable with values corresponding to the thresholds for the transitions from step 2.

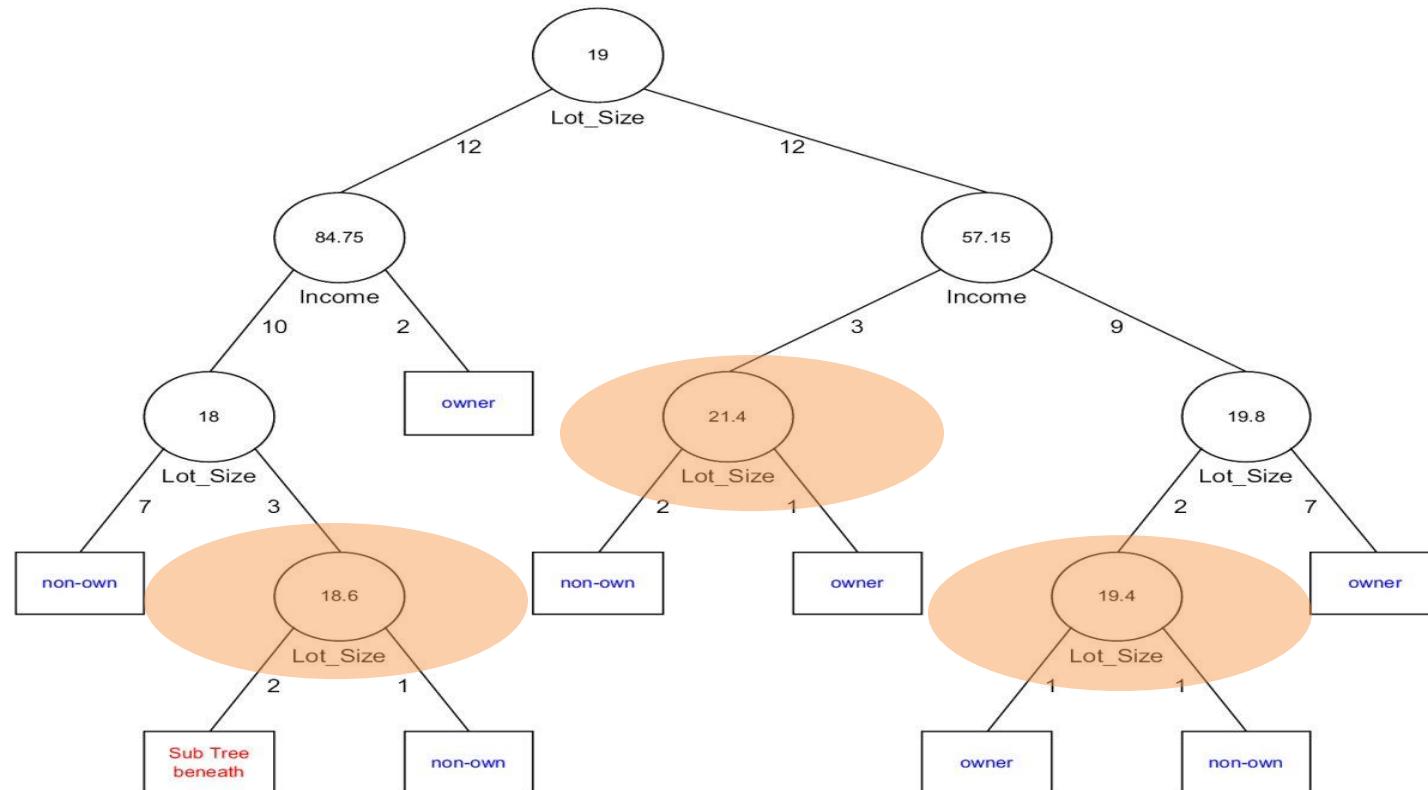
# Pruning



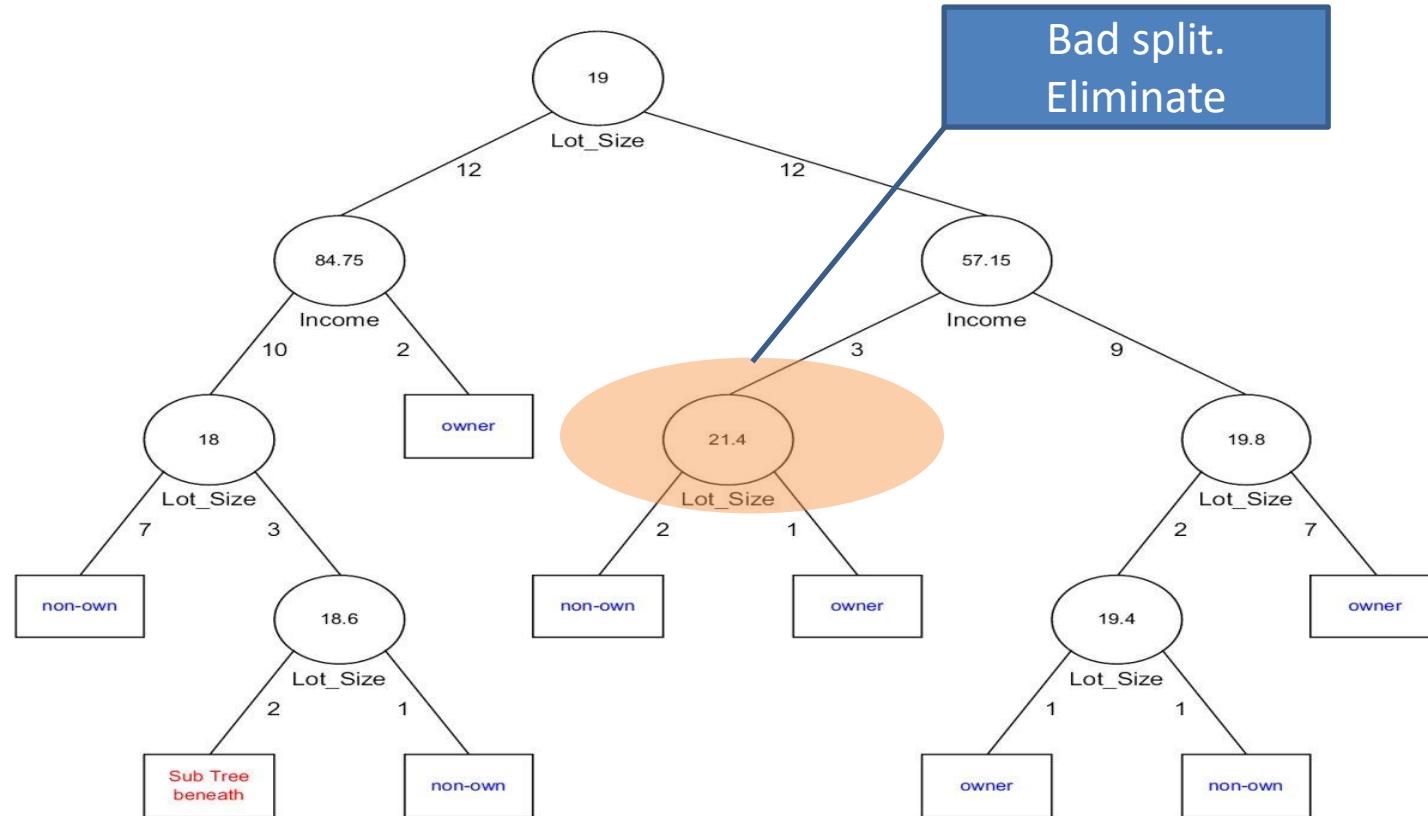
Big tree likely to over-fit

Eliminate splits that do not reduce error rate very much

# Splits Creating Small Groups

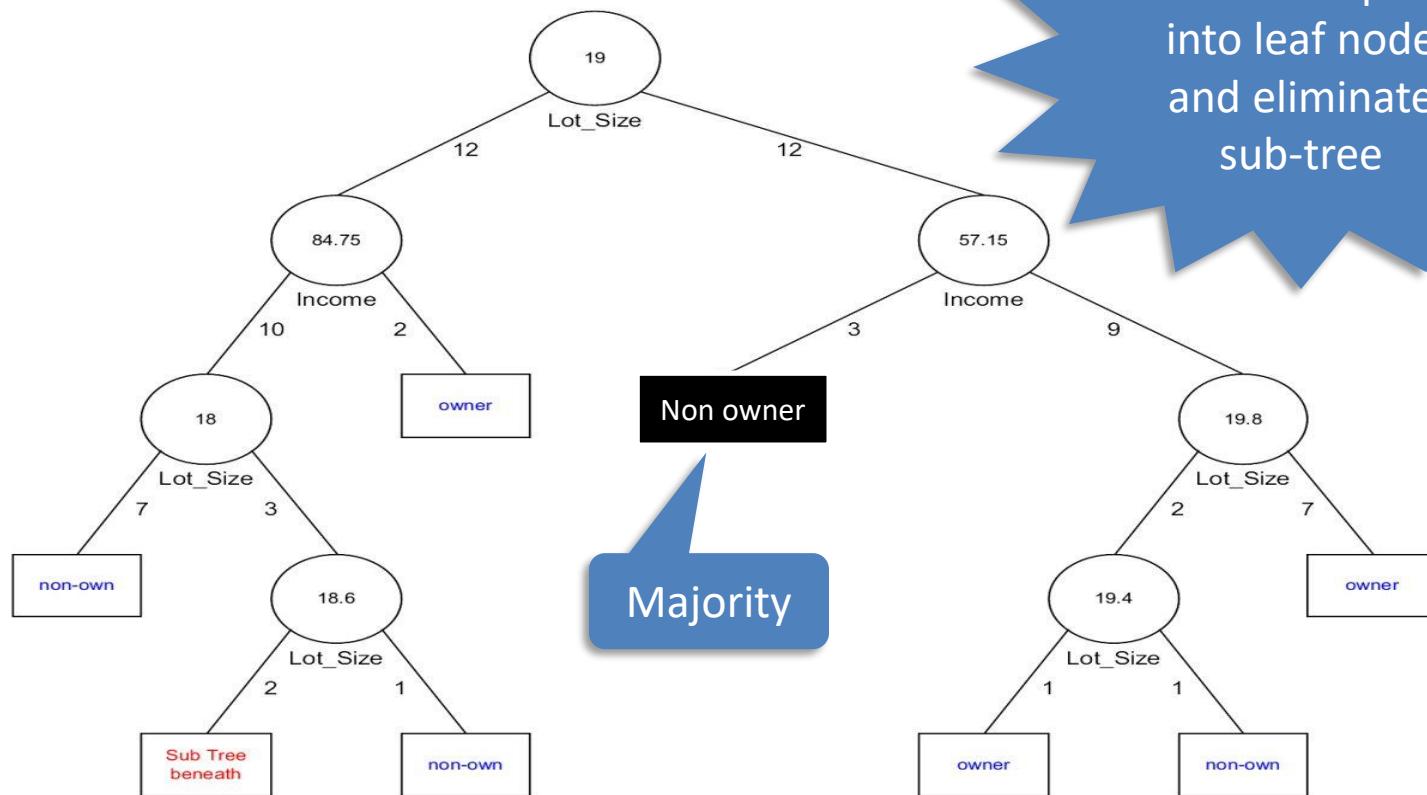


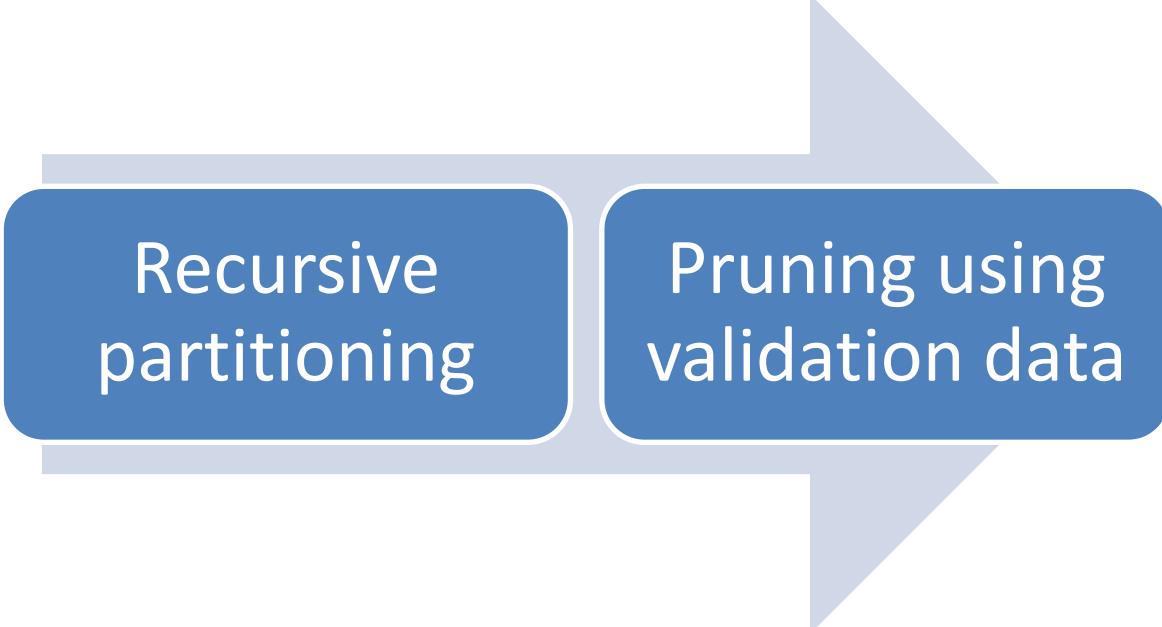
# Pruning Example



Bad split.  
Eliminate

# Pruning Example

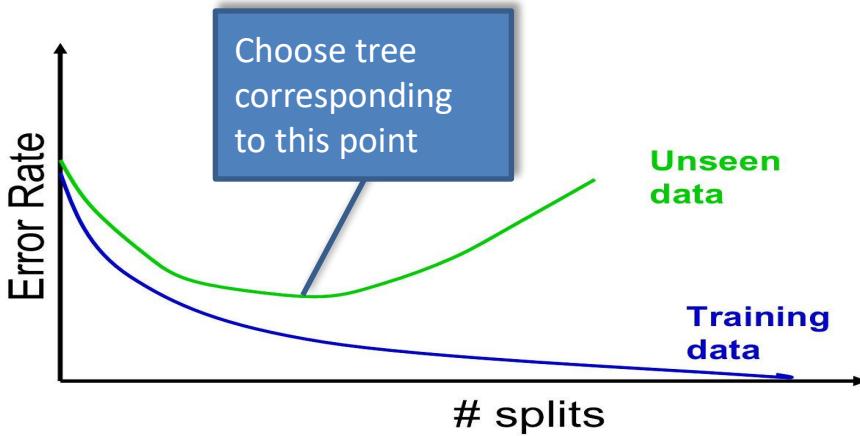




Recursive  
partitioning

Pruning using  
validation data

# Pruning Approach



Generate **best**  
trees of each  
size

Select tree with  
minimum **error**  
**rate on test**  
**data**

# Tree Implies Rules

## Advantages of rules

- Transparent
- Easy to explain
- Eg: Insurance -- able to explain why claim was denied, avoid lawsuits

## Advantages

## Disadvantages

Easy to use and interpret

Variable selection is automatic

Non-parametric

Missing data is not serious

Only horizontal and vertical splits

Sequential – could miss good splits



# Demo IRIS

WIKIPEDIA  
The Free Encyclopedia

Article

Talk



Not logged in Talk Contributions Create account Log in

Read

Edit

View history

Search Wikipedia



# Iris flower data set

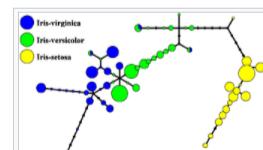
From Wikipedia, the free encyclopedia

[Main page](#)[Contents](#)[Current events](#)[Random article](#)[About Wikipedia](#)[Contact us](#)[Donate](#)[Contribute](#)[Help](#)[Learn to edit](#)[Community portal](#)[Recent changes](#)[Upload file](#)[Tools](#)[What links here](#)[Related changes](#)[Special pages](#)[Permanent link](#)[Page information](#)[Cite this page](#)[Wikidata item](#)[Print/export](#)[Download as PDF](#)[Printable version](#)[Languages](#)[العربية](#)[Español](#)[فارسی](#)[Français](#)[Italiano](#)[Русский](#)

## Contents [hide]

- 1 Use of the data set
- 2 Data set
  - 2.1 R code illustrating usage
  - 2.2 Python code illustrating usage
- 3 See also
- 4 References
- 5 External links

## Use of the data set [edit]

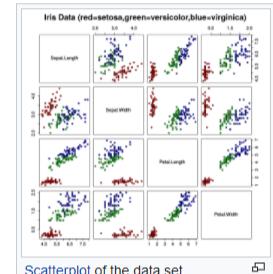


An example of the so-called "metro map" for the Iris data set.<sup>[4]</sup> Only a small fraction of Iris-virginica is mixed with Iris-versicolor. All other samples of the different Iris species belong to the different nodes.

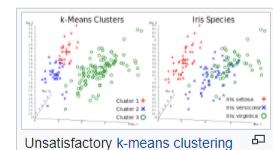
Based on Fisher's linear discriminant model, this data set became a typical test case for many statistical classification techniques in machine learning such as support vector machines.<sup>[5]</sup>

The use of this data set in cluster analysis however is not common, since the data set only contains two clusters with rather obvious separation. One of the clusters contains Iris setosa, while the other cluster contains both Iris virginica and Iris versicolor and is not separable without the species information Fisher used. This makes the data set a good example to explain the difference between supervised and unsupervised techniques in data mining: Fisher's linear discriminant model can only be obtained when the object species are known: class labels and clusters are not necessarily the same.<sup>[6]</sup>

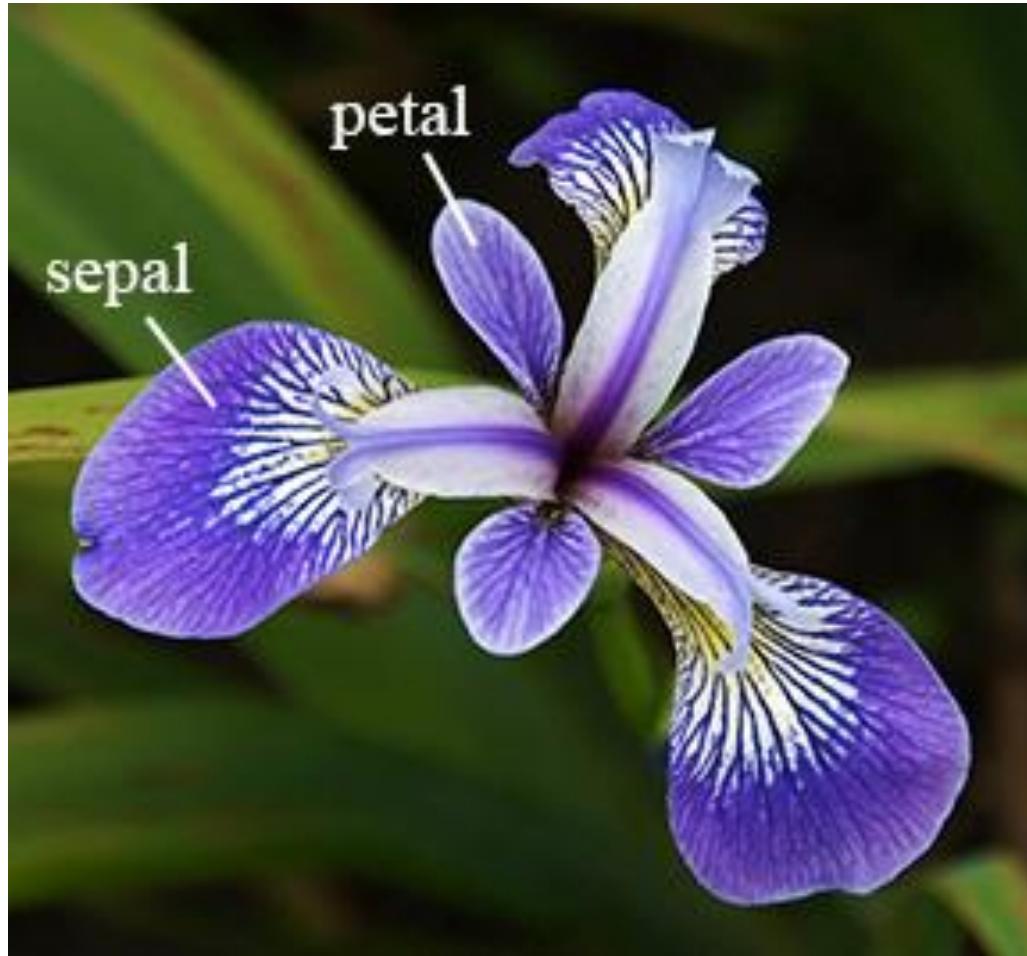
Nevertheless, all three species of Iris are separable in the projection on the nonlinear and branching principal component.<sup>[7]</sup> The data set is approximated by the closest tree with some penalty for the excessive number of nodes, bending and stretching. Then the so-called "metro map" is constructed.<sup>[4]</sup> The data points are projected into the closest node. For each node the pie diagram of the projected points is prepared. The area of the pie is proportional to the number of the projected points. It is clear from the diagram (left) that the absolute majority of the samples of the different Iris species belong to the different nodes. Only a small fraction of Iris virginica is mixed with Iris-versicolor (the mixed blue-green nodes in the diagram). Therefore, the three species of Iris (Iris setosa, Iris virginica and Iris versicolor) are separable by the



Scatterplot of the data set



Unsatisfactory K-means clustering (the data cannot be clustered into the known classes) and actual species visualized using ELKI



Pr... C in\_class\_coding.py

ITP449\_Fall2020 C

Class .idea Files In Class Coding in\_class\_coding.py

Run: in\_class\_coding ×

Number of rows: n\_samples

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa
..	...	...	...	...	...
145	6.7	3.0	5.2	2.3	virginica
146	6.3	2.5	5.0	1.9	virginica
147	6.5	3.0	5.2	2.0	virginica
148	6.2	3.4	5.4	2.3	virginica
149	5.9	3.0	5.1	1.8	virginica

Column: Feature

No of columns: n\_features

Row: Sample

1: Project Z: Structure 2: Favorites

4: Run 6: Problems Terminal Python Console

8:1 CRLF UTF-8 4 spaces Python 3.8 (Class)

Event Log

1: Project ITP449\_Fall2020 > Class > In Class Coding > in\_class\_coding.py 6 156

1 import seaborn as sns  
2  
3 iris = sns.load\_dataset('iris')  
4 print(iris)

Run: in\_class\_coding

Z: Structure 1: Structure 2: Favorites

Dependent variable: Target array **y** 1D with length n samples contained in array or Series

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa
..	...	...	...	...	...
145	6.7	3.0	5.2	2.3	virginica
146	6.3	2.5	5.0	1.9	virginica
147	6.5	3.0	5.2	2.0	virginica
148	6.2	3.4	5.4	2.3	virginica
149	5.9	3.0	5.1	1.8	virginica

4: Run 6: Problems Terminal Python Console Event Log

8:1 CRLF UTF-8 4 spaces Python 3.8 (Class)

## **Do the following:**

Import iris dataset from iris.csv and store as a dataframe

Print the dataframe

## 1: Project in\_class\_coding.py

```
1 import pandas as pd
2 import os
3
4 os.chdir('C:/Users/Reza/Desktop/ITP449_Fall2020/Class/Files')
5 pd.set_option('display.max_columns', None)
6
7 irisDf = pd.read_csv('iris.csv')
8
9 print(irisDf.head())
10
```

⚠ 3 ✅ 258

## 2: Structure

## 3: Favorites

## 4: Run

## Run: in\_class\_coding



```
C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Class/In Class Coding/in_class_coding.py"
   Sepal length  Sepal width  Petal length  Petal width    Species
0            5.1         3.5          1.4         0.2  I. setosa
1            4.9         3.0          1.4         0.2  I. setosa
2            4.7         3.2          1.3         0.2  I. setosa
3            4.6         3.1          1.5         0.2  I. setosa
4            5.0         3.6          1.4         0.3  I. setosa
```

```
Process finished with exit code 0
```

**Do the following:**

Create feature matrix and response vector: X, y

## 1: Project in\_class\_coding.py

```
1 import pandas as pd
2 import os
3
4 os.chdir('C:/Users/Reza/Desktop/ITP449_Fall2020/Class/Files')
5 pd.set_option('display.max_columns', None)
6
7 irisDf = pd.read_csv('iris.csv')
8
9 X = irisDf[['Sepal length', 'Sepal width', 'Petal length', 'Petal width']]
10 y = irisDf['Species']
11
```

3 258

## 2: Structure

## 3: Favorites

## Run: in\_class\_coding



```
▶ C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Class/In Class Co
```

Process finished with exit code 0

## 4: Favorites

## 5: Run

```
4: Run TODO 6: Problems Terminal Python Console Event Log
```

## **Do the following:**

Partition the X and y into train and test partitions - 70/30.



in\_class\_coding.py

```
1 import pandas as pd
2 import os
3
4 os.chdir('C:/Users/Reza/Desktop/ITP449_Fall2020/Class/Files')
5 pd.set_option('display.max_columns', None)
6
7 irisDf = pd.read_csv('iris.csv')
8
9 X = irisDf[['Sepal length', 'Sepal width', 'Petal length', 'Petal width']]
10 y = irisDf['Species']
11
12 from sklearn.model_selection import train_test_split
13
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=2020, stratify=y)
```

Run:  in\_class\_coding ×

—

C:\Users\Reza\Desktop\ITP449\_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449\_Fall2020/Class/In Class Code/Assignment 1.py"



**Do the following:**

Train the decision tree classifier

## Project 1: in\_class\_coding.py

```
1 import pandas as pd
2 import os
3
4 os.chdir('C:/Users/Reza/Desktop/ITP449_Fall2020/Class/Files')
5 pd.set_option('display.max_columns', None)
6
7 irisDf = pd.read_csv('iris.csv')
8
9 X = irisDf[['Sepal length', 'Sepal width', 'Petal length', 'Petal width']]
10 y = irisDf['Species']
11
12 from sklearn.model_selection import train_test_split
13
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=2020, stratify=y)
15
16 from sklearn.tree import DecisionTreeClassifier
17
18 dt = DecisionTreeClassifier(criterion='entropy')
19 dt.fit(X_train, y_train)
```

criterion: *gini* or *entropy*

## 2: Structure

## Favorites

Run: in\_class\_coding



&gt;&gt;

&gt;&gt;

▶ 4: Run

TODO

6: Problems

Terminal

Python Console

Event Log

## **Do the following:**

- Predict the species using the test partition.
- Plot the confusion matrix
- Print the accuracy

## Project 1: in\_class\_coding.py

```
18 dt = DecisionTreeClassifier(criterion='entropy')
19 dt.fit(X_train, y_train)
20
21 y_pred = dt.predict(X_test)
22
23 from sklearn import metrics
24
25 cf = metrics.confusion_matrix(y_test, y_pred)
26 print(cf)
27
28 print('Accuracy =', metrics.accuracy_score(y_test, y_pred))
29
30 import matplotlib.pyplot as plt
31
32 metrics.plot_confusion_matrix(dt, X_test, y_test)
33 plt.show()
```

⚠ 7 ✅ 259

## Structure

## Favorites

## Run: in\_class\_coding



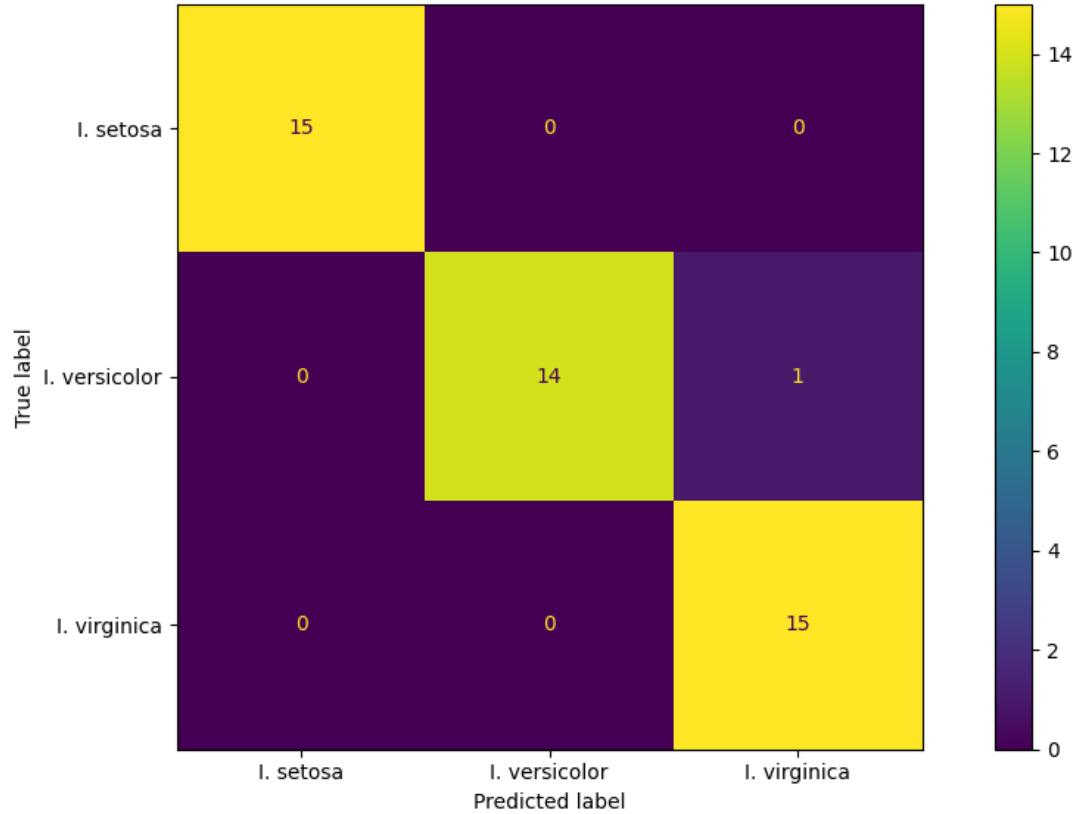
```
[[15  0  0]
 [ 0 14  1]
 [ 0  0 15]]
```

Accuracy = 0.9777777777777777



▶ 4: Run TODO 6: Problems Terminal Python Console

Event Log



**Do the following:**

Display the decision tree

## Project 1: Structure Z: in\_class\_coding.py

```
16  from sklearn.tree import DecisionTreeClassifier
17
18  dt = DecisionTreeClassifier(criterion='entropy')
19  dt.fit(X_train, y_train)
20
21  y_pred = dt.predict(X_test)
22
23  from sklearn import metrics
24
25  cf = metrics.confusion_matrix(y_test, y_pred)
26  print(cf)
27
28  print('Accuracy =', metrics.accuracy_score(y_test, y_pred))
29
30  import matplotlib.pyplot as plt
31
32  metrics.plot_confusion_matrix(dt, X_test, y_test)
33  plt.show()
34
```

A 8 ✓ 259

## Favorites

## Run: in\_class\_coding



&gt;&gt;

▶ 4: Run

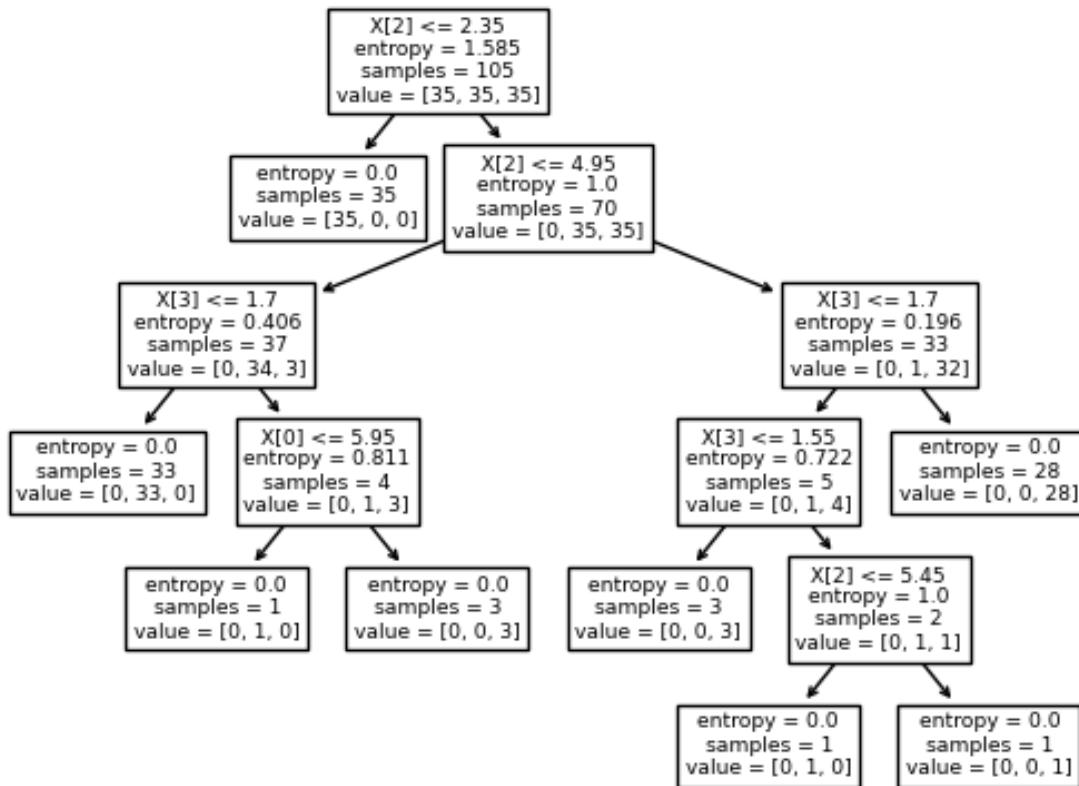
Todo

Problems

Terminal

Python Console

Event Log



## **Do the following:**

Format the tree to be more descriptive and save the decision tree as a png

## Project 1: in\_class\_coding.py

```
23 from sklearn import metrics
24
25 cf = metrics.confusion_matrix(y_test, y_pred)
26 print(cf)
27
28 print('Accuracy =', metrics.accuracy_score(y_test, y_pred))
29
30 import matplotlib.pyplot as plt
31
32 metrics.plot_confusion_matrix(dt, X_test, y_test)
33 plt.show()
34
35 from sklearn import tree
36
37 plt.figure(2)
38 fn = X.columns
39 cn = y.unique()
40 irisTree = tree.plot_tree(dt, feature_names=fn, class_names=cn, filled=True)
41 plt.savefig('IrisDT.png')
42 plt.show()
43
```

⚠ 9 ✅ 264

## Structure

## Favorites

## Run: in\_class\_coding



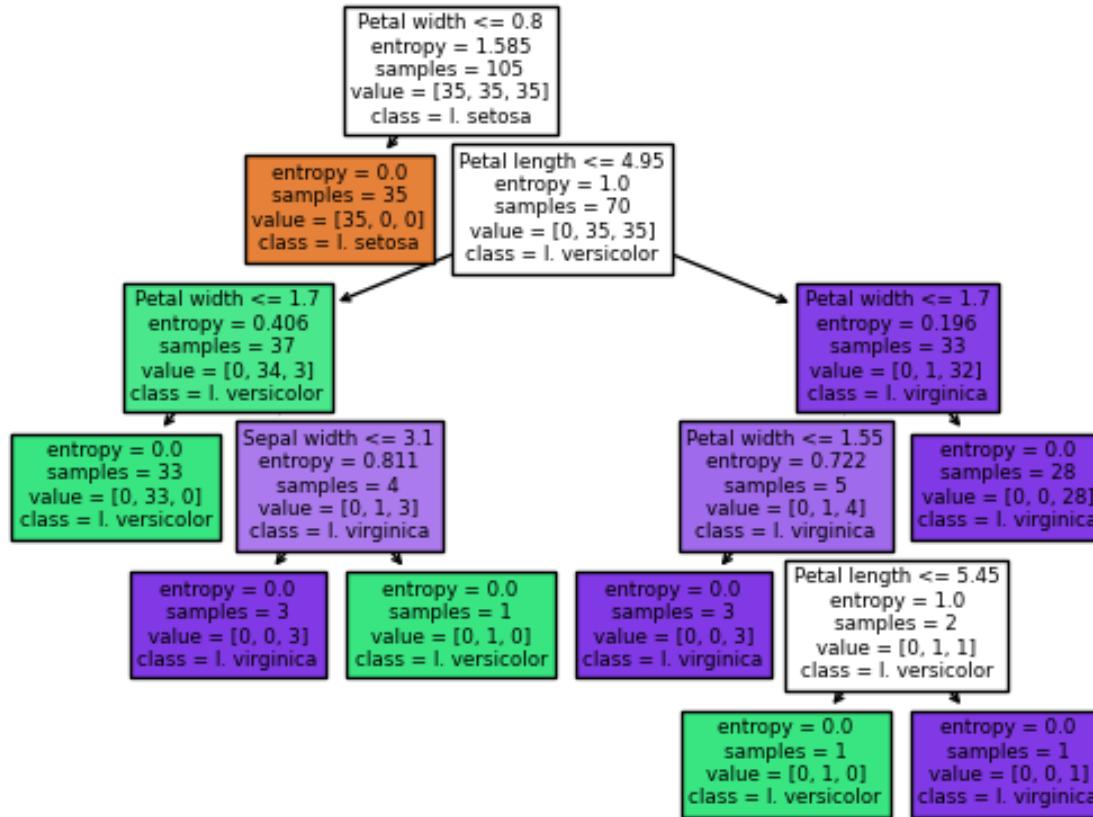
▶ 4: Run

Todo 6: Problems

Terminal

Python Console

Event Log



**Do the following:**

Show Feature importance

## Project 1: in\_class\_coding.py

```
34
35     from sklearn import tree
36
37     plt.figure(2)
38     fn = X.columns
39     cn = y.unique()
40     irisTree = tree.plot_tree(dt, feature_names=fn, class_names=cn, filled=True)
41     plt.savefig('IrisDT.png')
42     plt.show()
43
44     # Feature Importance
45
46     print('Feature Importance:', dt.feature_importances_)
47
```

⚠ 9 ✅ 264 ⏪

## Z: Structure

## ★ 2: Favorites

## Run: in\_class\_coding

```
C:\Users\Reza\Desktop\ITP449_Fall2020\Class\venv\Scripts\python.exe "C:/Users/Reza/Desktop/ITP449_Fall2020/Class/In Class Coding/in_class_coding.py"
[[15  0  0]
 [ 0 14  1]
 [ 0  0 15]]
Accuracy = 0.9777777777777777
Feature Importance: [0.0096721  0.01949941  0.87089281  0.09993569]
```

Petal Length is the most important feature.