



# Decision Trees: Interpretable Classification (and Regression)

Week 2

January 29, 2024

Spring 24 | CUSP-GX 7033 – Machine Learning for Cities | Dr. Anton Rozhkov

# Today's Outline

- The prediction problem (classification and regression)
- Overview of classification methods; discussion of interpretable vs. black-box predictors; feature importance
- Interpretable prediction: Rule-based, instance-based, and model-based
- Rule-based Learning: Decision trees for classification and regression

# Overview of Classification Methods

# The Prediction Problem

Most machine learning tasks are about making predictions.

This means that given a number of training examples, the system needs to be able to make good predictions for (generalize to) examples it has never seen before.

Having a good performance measure on the training data is good but insufficient; the true goal is to perform well on new instances.

# Classification vs Regression

**Classification** predictive modeling is the task of approximating a mapping function ( $f$ ) from input variables ( $X$ ) to discrete output variables ( $y$ ).

**Regression** predictive modeling is the task of approximating a mapping function ( $f$ ) from input variables ( $X$ ) to a continuous output variable ( $y$ ).

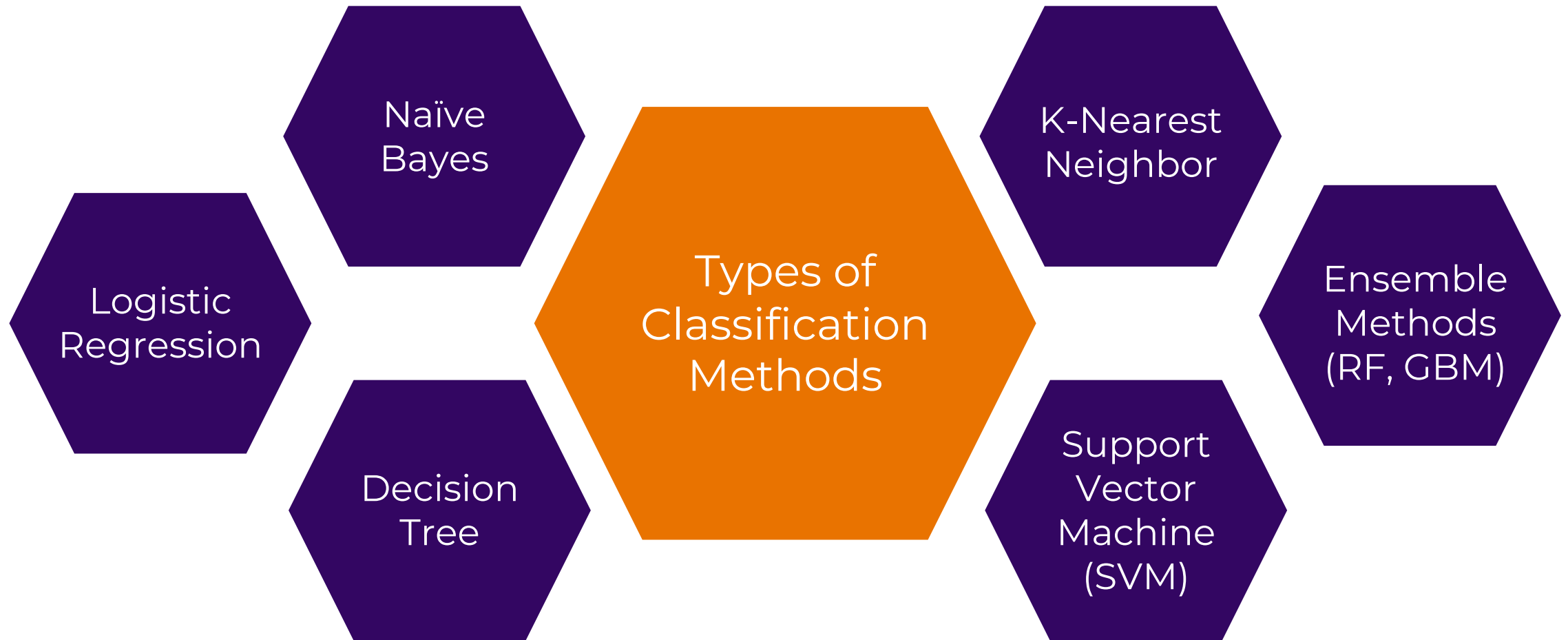
# Classification vs Regression

**Classification** predictive modeling is the task of approximating a mapping function ( $f$ ) from input variables ( $X$ ) to discrete output variables ( $y$ ).



**Regression** predictive modeling is the task of approximating a mapping function ( $f$ ) from input variables ( $X$ ) to a continuous output variable ( $y$ ).

# Classification Methods



# Evaluating Classification Methods

- **Accuracy**

- classifier accuracy: predicting the class label
- regression accuracy: guessing the value of predicted attributes

- **Speed**

- time to construct the model (training time)
- time to use the model (classification/prediction time)

- **Robustness**: handling noise and missing values

- **Scalability**: efficiency on large databases

- **Interpretability**

- understanding and insight provided by the model

- **Other measures**, e.g., the goodness of rules, such as decision tree size or compactness of classification rules



# Interpretable Prediction

# Idea of a Black Box



# Introductions

Logistic regression  
Decision trees  
...

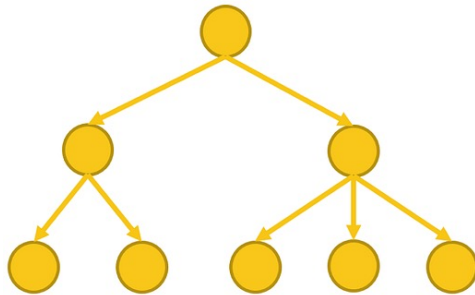
...

Random forest

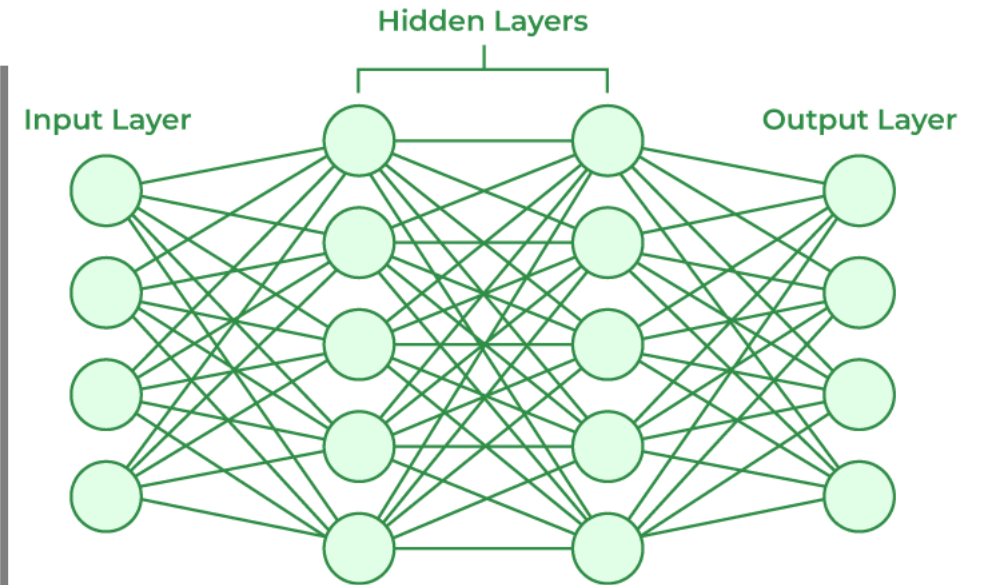
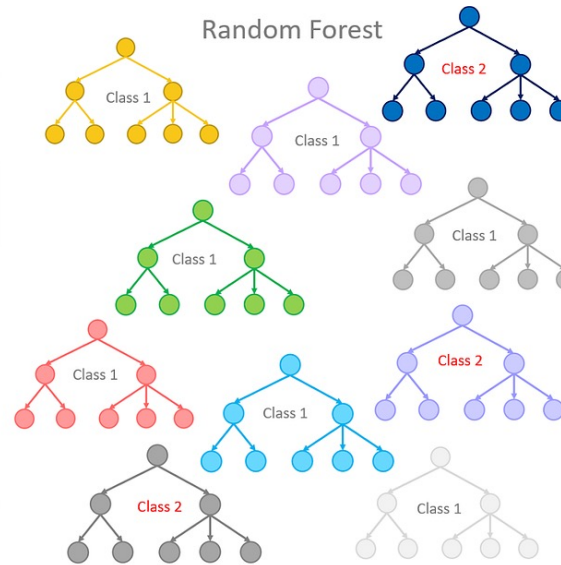
...

Artificial Neural Networks

Single Decision Tree



**Interpretability  
(transparent)**



**“Black box”  
(opaque)**

# Idea of a Black Box



## To Explain or to Predict?

Galit Shmueli

**Abstract.** Statistical modeling is a powerful tool for developing and testing theories by way of causal explanation, prediction, and description. In many disciplines there is near-exclusive use of statistical modeling for causal explanation and the assumption that models with high explanatory power are inherently of high predictive power. Conflation between explanation and prediction is common, yet the distinction must be understood for progressing scientific knowledge. While this distinction has been recognized in the philosophy of science, the statistical literature lacks a thorough discussion of the many differences that arise in the process of modeling for an explanatory versus a predictive goal. The purpose of this article is to clarify the distinction between explanatory and predictive modeling, to discuss its sources, and to reveal the practical implications of the distinction to each step in the modeling process.

**Key words and phrases:** Explanatory modeling, causality, predictive modeling, predictive power, statistical strategy, data mining, scientific research.

### 1. INTRODUCTION

Looking at how statistical models are used in different scientific disciplines for the purpose of theory building and testing, one finds a range of perceptions regarding the relationship between causal explanation and empirical prediction. In many scientific fields such as economics, psychology, education, and environmental science, statistical models are used almost exclusively for causal explanation, and models that possess high explanatory power are often assumed to inherently possess predictive power. In fields such as natural language processing and bioinformatics, the focus is on empirical prediction with only a slight and indirect relation to causal explanation. And yet in other research fields, such as epidemiology, the emphasis on causal explanation versus empirical prediction is more mixed. Statistical modeling for description, where the purpose is to capture the data structure parsimoniously, and which is the most commonly developed within the field of statistics, is not commonly used for theory building and testing in other disciplines. Hence, in this article I

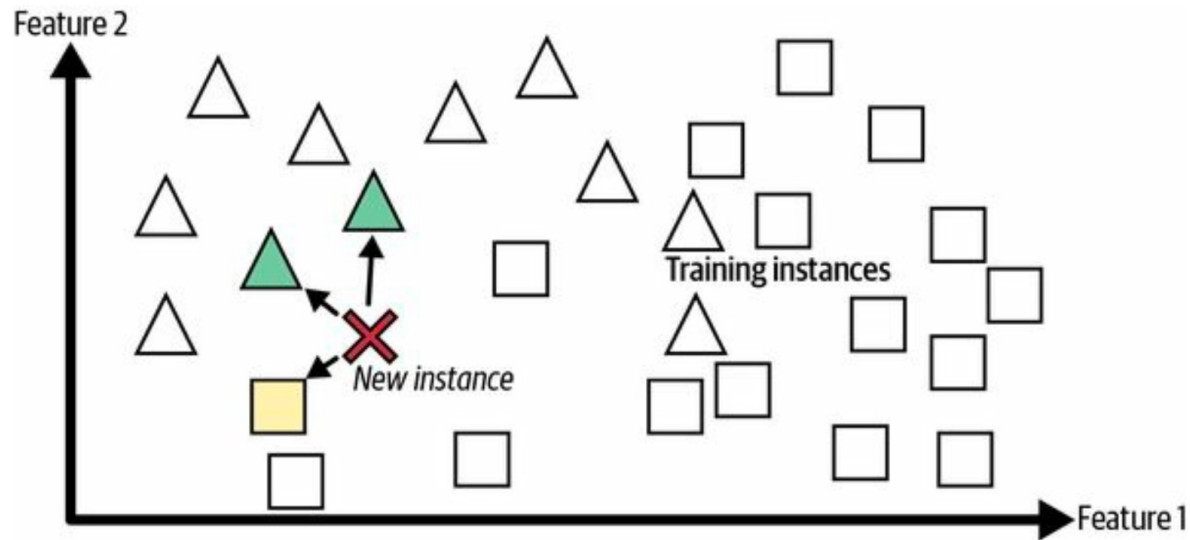
focus on the use of statistical modeling for causal explanation and for prediction. My main premise is that the two are often conflated, yet the causal versus predictive distinction has a large impact on each step of the statistical modeling process and on its consequences. Although not explicitly stated in the statistics methodology literature, applied statisticians instinctively sense that predicting and explaining are different. This article aims to fill a critical void: to tackle the distinction between explanatory modeling and predictive modeling.

Clearing the current ambiguity between the two is critical not only for proper statistical modeling, but more importantly, for proper scientific usage. Both explanation and prediction are necessary for generating and testing theories, yet each plays a different role in doing so. The lack of a clear distinction within statistics has created a lack of understanding in many disciplines of the difference between building sound explanatory models versus creating powerful predictive models, as well as confusing explanatory power with predictive power. The implications of this omission and the lack of clear guidelines on how to model for explanatory versus predictive goals are considerable for both scientific research and practice and have also contributed to the gap between academia and practice.

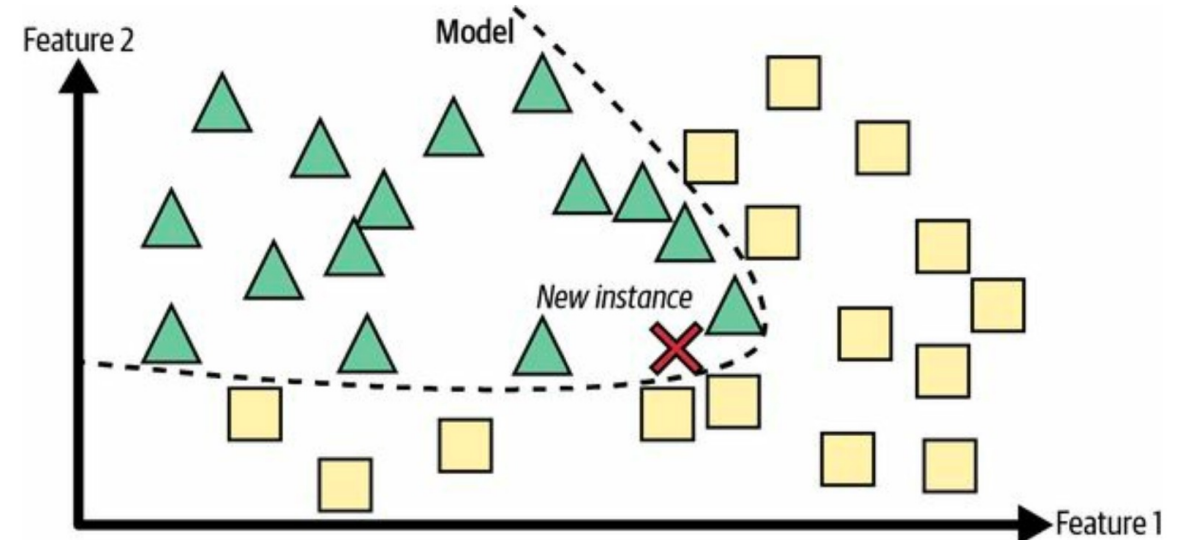
I start by defining what I term *explaining* and *predicting*. These definitions are chosen to reflect the dis-

*Galit Shmueli is Associate Professor of Statistics, Department of Decision, Operations and Information Technologies, Robert H. Smith School of Business, University of Maryland, College Park, Maryland 20742, USA (e-mail: gshmueli@umd.edu).*

# Instance-based and Model-based Learning



**Instance-based**



**Model-based**

# Instance-based and Model-based Learning

## Instance-based

## Model-based

### Generalization

Simply memorize the training examples

Learn a generalizable model

### Scalability

Slow and memory-intensive with large datasets

More scalable without storing the training data

### Interpretability

Effective for small or medium-sized datasets

Good for all datasets

# Rule-based Learning

Rule-based systems are computer programs utilizing *if-then* rules for decision-making and task execution.

## Strengths:

- Structured Decision Making: offers a clear and structured approach to problem-solving.
- Transparency: the logic behind decisions is explicit, providing transparency.

## Challenges:

- Optimal Solutions: rule-based systems may not always find the most optimal solution.
- Dependency on Rule Quality: effectiveness is tied to the accuracy and comprehensiveness of rules.

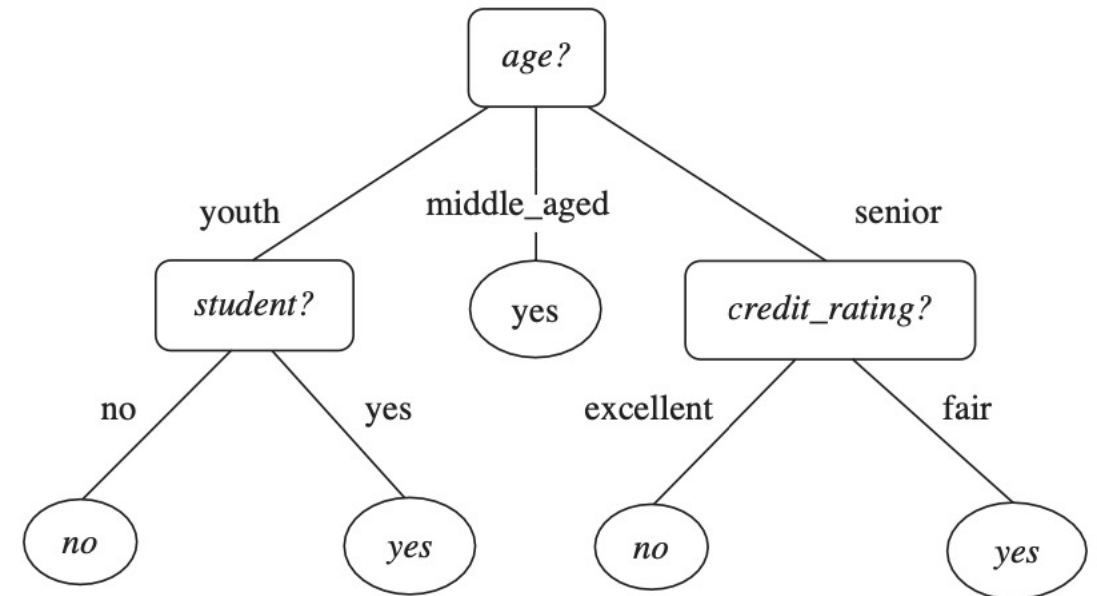
# Decision Trees



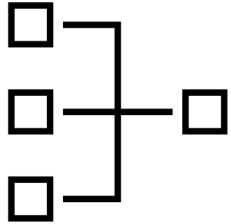
# Decision Trees

**Decision Trees** (DTs) are a non-parametric supervised learning method used for classification and regression.

- The simplest method to find rules.
- Provide “if condition, then consequence” rules.
- Simple to understand and to interpret. Trees can be visualized.
- Requires little data preparation.

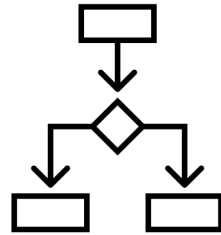


# History



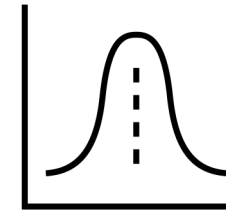
*Concept learning systems*

by E. B. Hunt, J. Marin,  
and P. T. Stone



Development of a  
decision tree algorithm  
known as *ID3 (Iterative  
Dichotomiser)* + C4.5

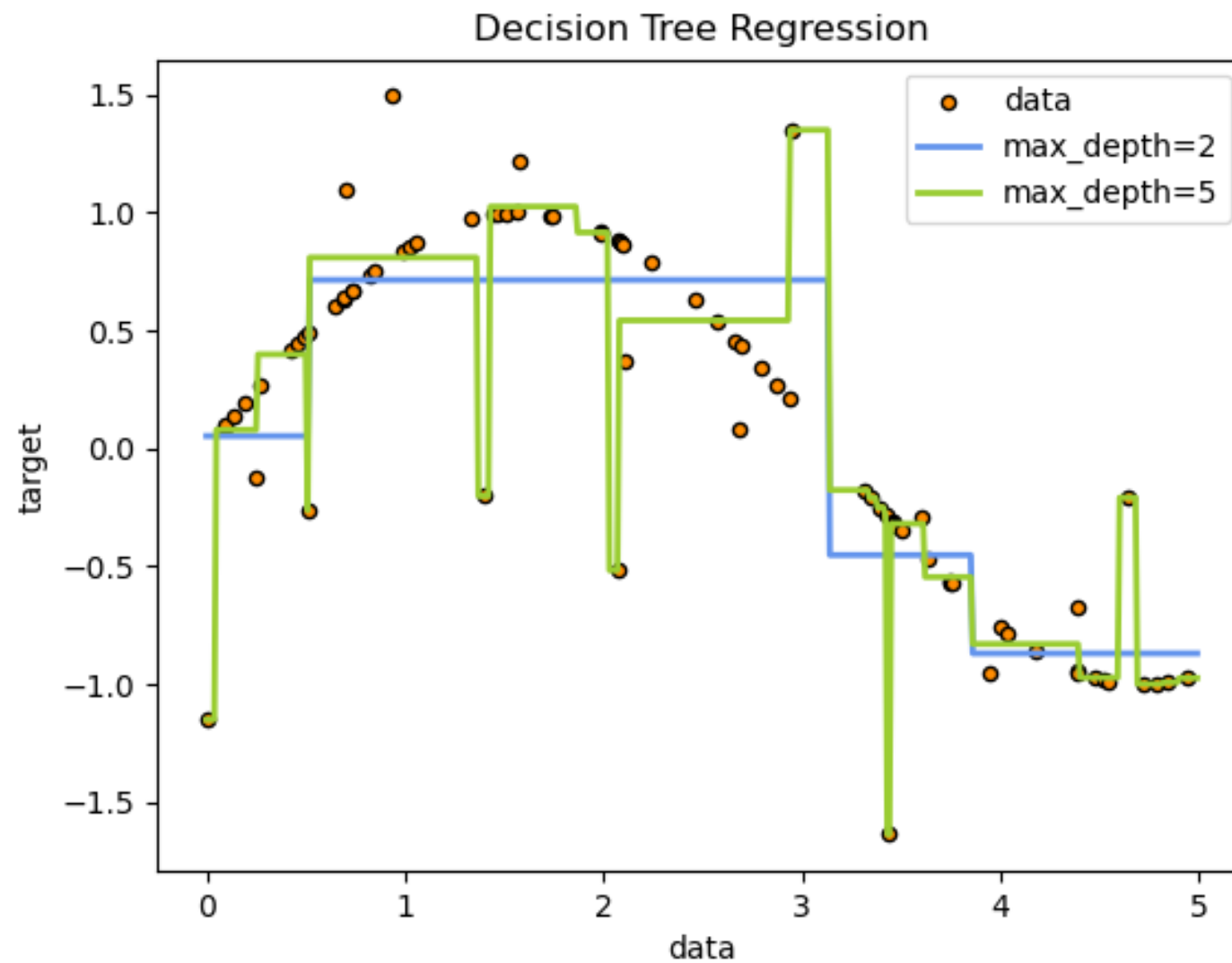
By J. Ross Quinlan



*Book: Classification and  
Regression Trees (CART),  
which described the  
generation of binary  
decision trees*

by L. Breiman, J. Friedman,  
R. Olshen, and C. Stone

# Example



# When to Use Decision Trees? (1/2)

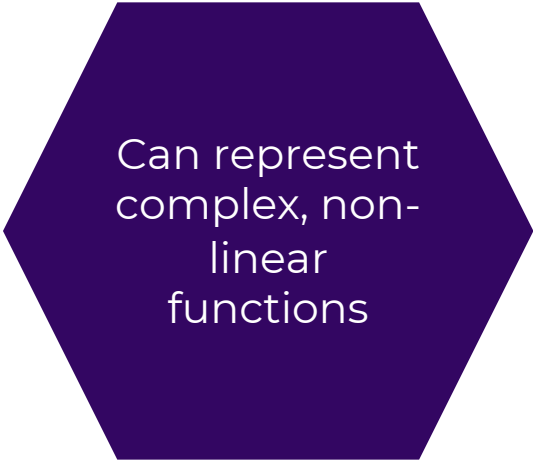
We have a dataset, with some attribute we want to predict

- We can do either classification or regression.
- Datasets with lots of attributes and/or lots of records are okay. But if there are lots of attributes and not that many records, you should probably use feature selection first to avoid overfitting.
- Datasets with discrete or real values, or both, are okay.

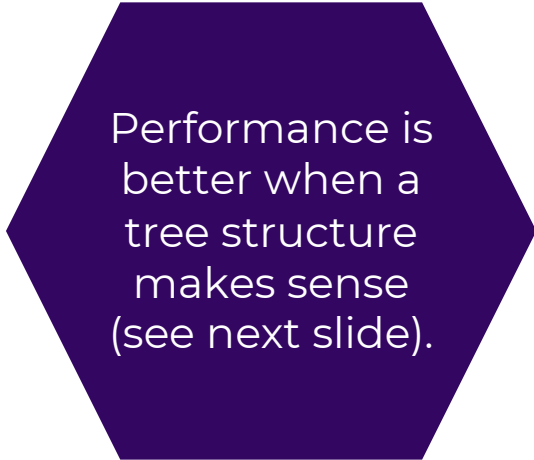
We want to be able to explain our prediction using a simple and interpretable set of rules

- We want to distinguish relevant from irrelevant attributes.
- We want to provide a set of decision steps (e.g., an “expert system” for medical diagnosis; we don’t want to perform irrelevant tests).
- If all we care about is prediction accuracy and not interpretability, we might want to use some “black box” classifier instead (e.g., neural network, support vector machine).

# When to Use Decision Trees? (2/2)



Can represent  
complex, non-  
linear  
functions

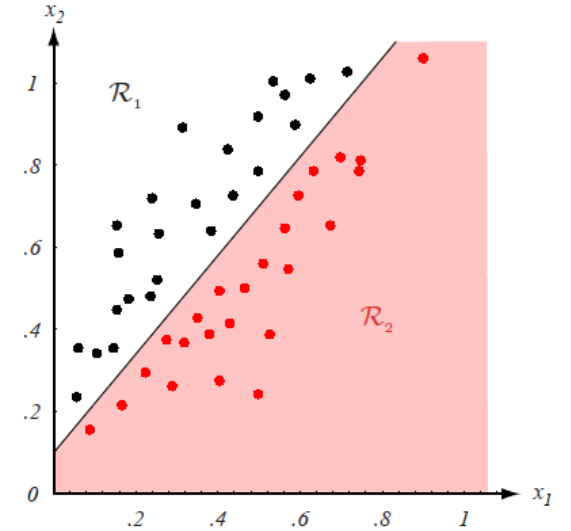


Performance is  
better when a  
tree structure  
makes sense  
(see next slide).

# Inductive Bias

**Inductive bias** describes how a prediction method generalizes to previously unseen examples.

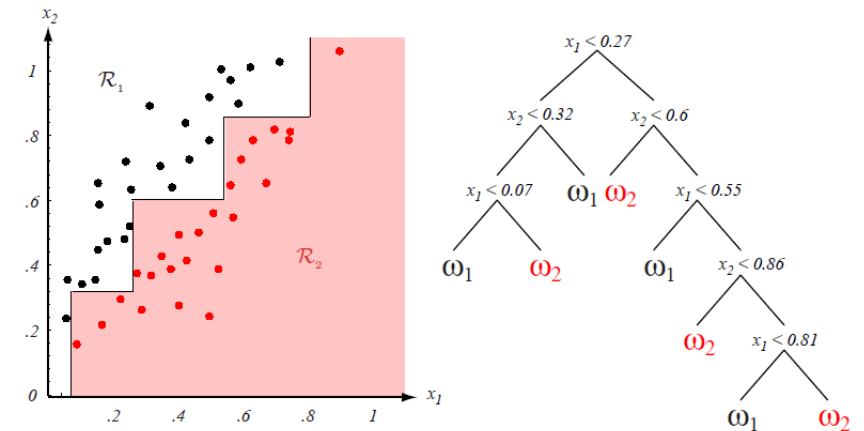
For example: binary logistic regression assumes a linear decision boundary between  $R_1(y = 1) > 0.5$  and  $R_2(y = 1) < 0.5$ .



What are the inductive biases of a decision tree?

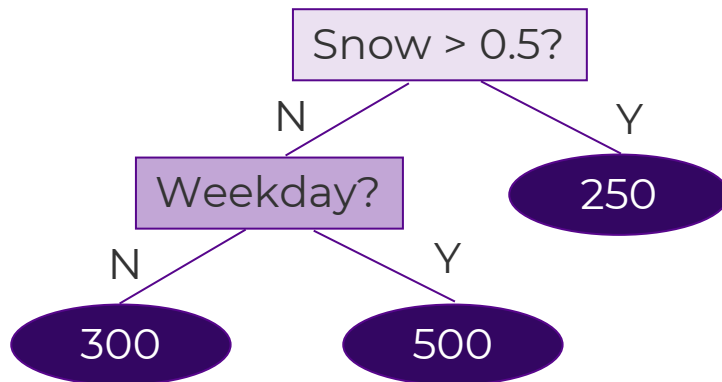
- Flexible – any function can be represented.
- Prefers shorter trees to longer trees.
- Prefers attributes with higher information gain.
- Splits on one attribute at a time.
- Splits are axis-aligned.

Suboptimal for linear functions

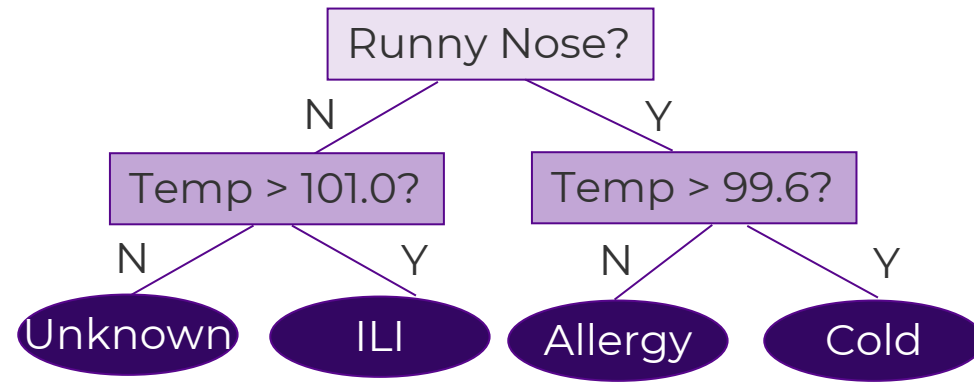


# Rule-based Learning with Decision Trees

- A decision tree is a set of rules that can be learned from data and used to predict an unknown value.
  - Unknown real value (regression): What is the expected incidence of car thefts in NYC on a given day?
  - Unknown category value (classification): What type of illness does patient X have, given their symptoms and demographic data?



How many thefts on Tuesday,  
January 3 (0.2 inches of snow)?



What do we predict for a patient  
with Temp = 100 and a runny nose?

# Learning Binary Decision Trees

- Step 1: Start with all data points in a single node. Predict the most common value (classification) or mean value (regression).

**BAD**

(5 good, 15 bad)

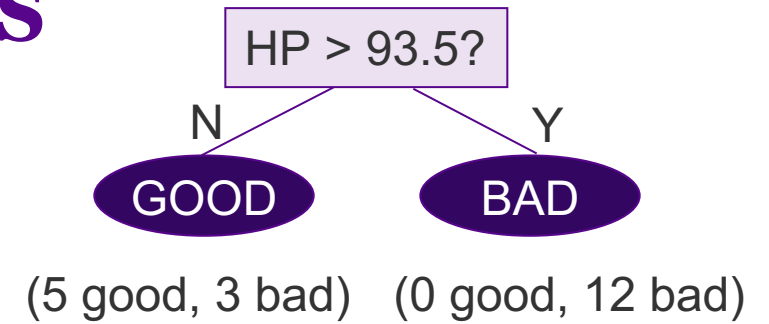
MPG, cylinders, HP, weight

good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
bad, 4, 81, light  
bad, 6, 95, medium  
good, 4, 93, light



# Learning Binary Decision Trees

- Step 1: Start with all data points in a single node. Predict the most common value (classification) or mean value (regression).
- Step 2: Choose the “best” binary decision rule and use it to split the data into two groups.

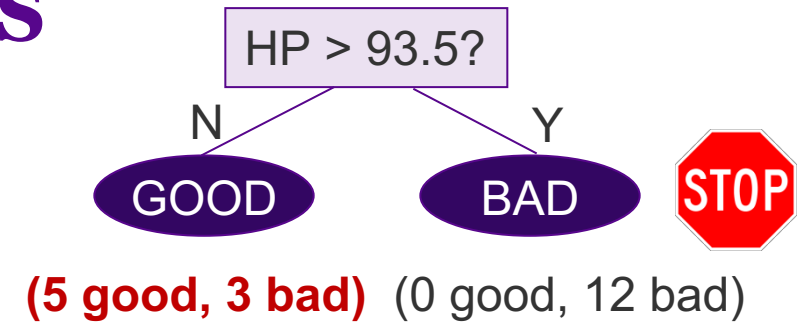


MPG, cylinders, HP, weight

good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
bad, 4, 81, light  
bad, 6, 95, medium  
good, 4, 93, light

# Learning Binary Decision Trees

- Step 1: Start with all data points in a single node. Predict the most common value (classification) or mean value (regression).
- Step 2: Choose the “best” binary decision rule and use it to split the data into two groups.
- Step 3: Repeat step 2 on each group until some stopping criterion is reached.



MPG, cylinders, HP, weight

**good, 4, 75, light**

**bad, 6, 90, medium**

bad, 4, 110, medium

bad, 8, 175, weighty

bad, 6, 95, medium

bad, 4, 94, light

bad, 4, 95, light

bad, 8, 139, weighty

bad, 8, 190, weighty

bad, 8, 145, weighty

bad, 6, 100, medium

**good, 4, 92, medium**

bad, 6, 100, weighty

bad, 8, 170, weighty

**good, 4, 89, medium**

**good, 4, 65, light**

**bad, 6, 85, medium**

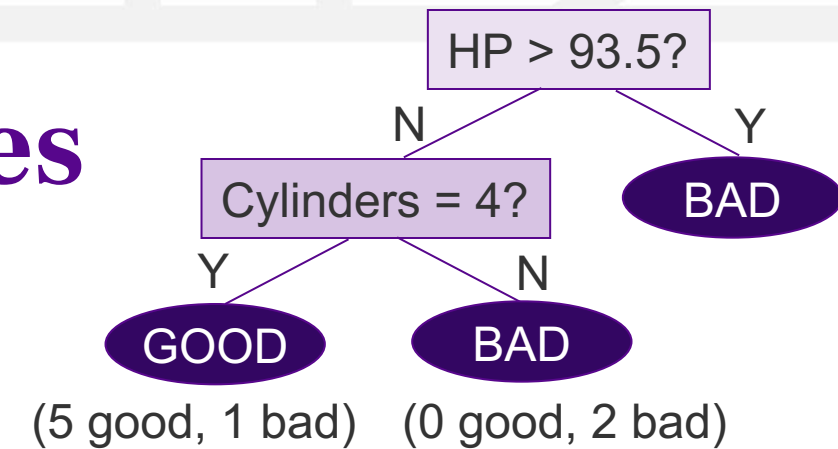
**bad, 4, 81, light**

bad, 6, 95, medium

**good, 4, 93, light**

# Learning Binary Decision Trees

- Step 1: Start with all data points in a single node. Predict the most common value (classification) or mean value (regression).
- Step 2: Choose the “best” binary decision rule and use it to split the data into two groups.
- Step 3: Repeat step 2 on each group until some stopping criterion is reached.



MPG, cylinders, HP, weight

**good, 4, 75, light**

**bad, 6, 90, medium**

bad, 4, 110, medium

bad, 8, 175, weighty

bad, 6, 95, medium

bad, 4, 94, light

bad, 4, 95, light

bad, 8, 139, weighty

bad, 8, 190, weighty

bad, 8, 145, weighty

bad, 6, 100, medium

**good, 4, 92, medium**

bad, 6, 100, weighty

bad, 8, 170, weighty

**good, 4, 89, medium**

**good, 4, 65, light**

**bad, 6, 85, medium**

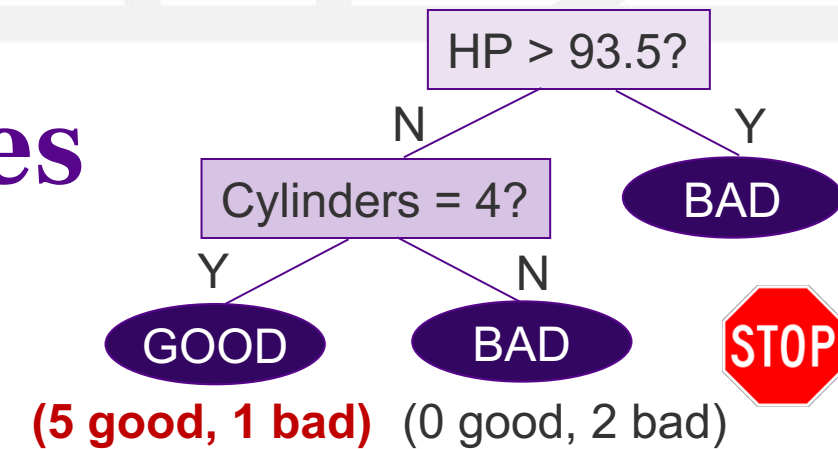
**bad, 4, 81, light**

bad, 6, 95, medium

**good, 4, 93, light**

# Learning Binary Decision Trees

- Step 1: Start with all data points in a single node. Predict the most common value (classification) or mean value (regression).
- Step 2: Choose the “best” binary decision rule and use it to split the data into two groups.
- Step 3: Repeat step 2 on each group until some stopping criterion is reached.



MPG, cylinders, HP, weight

**good, 4, 75, light**

bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium

**good, 4, 92, medium**

bad, 6, 100, weighty  
bad, 8, 170, weighty

**good, 4, 89, medium**

**good, 4, 65, light**

bad, 6, 85, medium

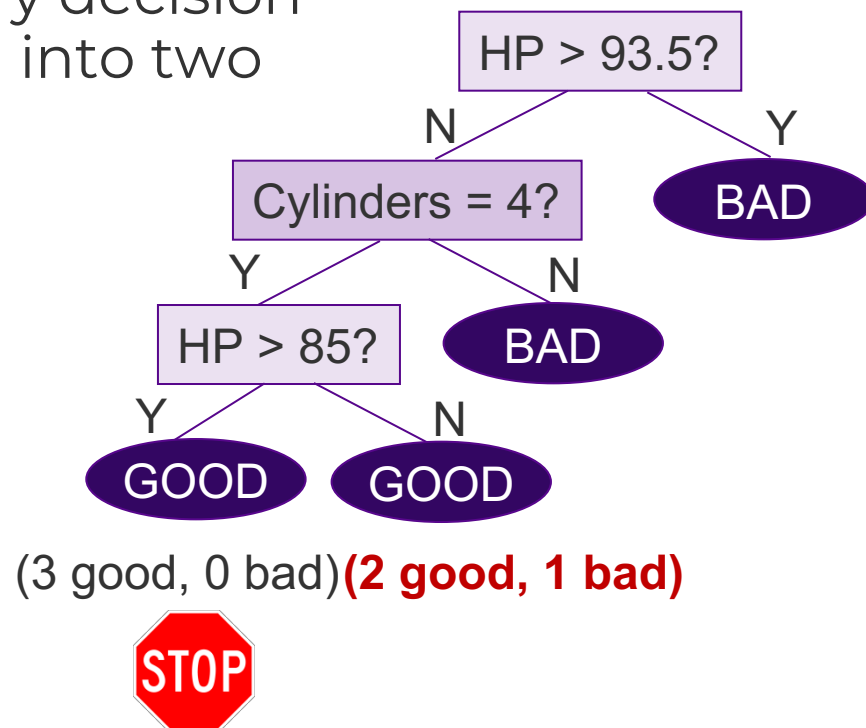
**bad, 4, 81, light**

bad, 6, 95, medium

**good, 4, 93, light**

# Learning Binary Decision Trees

- Step 1: Start with all data points in a single node. Predict the most common value (classification) or mean value (regression).
- Step 2: Choose the “best” binary decision rule and use it to split the data into two groups.
- Step 3: Repeat step 2 on each group until some stopping criterion is reached.



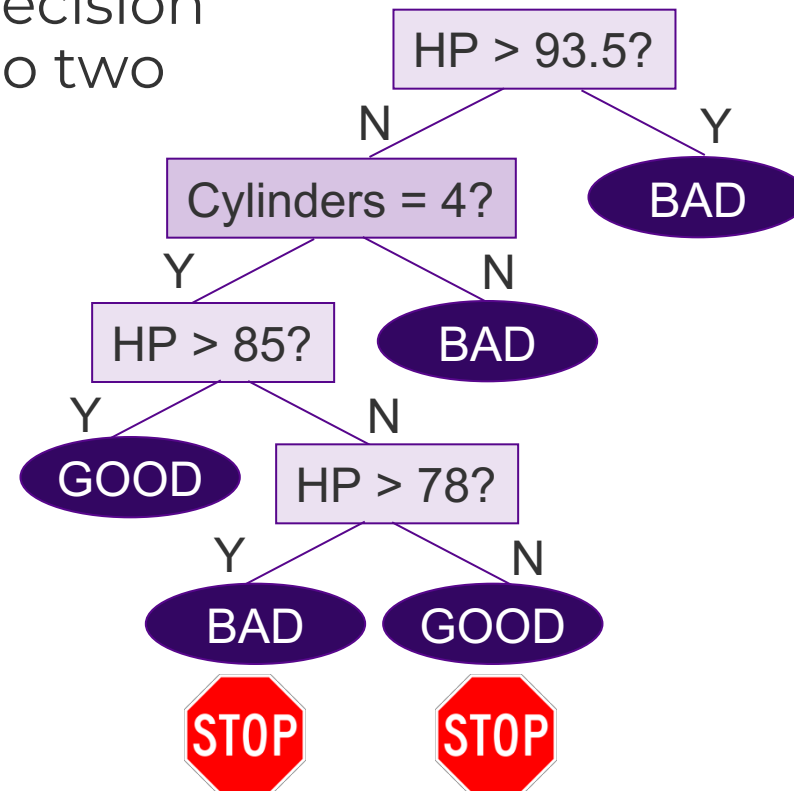
MPG, cylinders, HP, weight

**good, 4, 75, light**

bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
**good, 4, 65, light**  
bad, 6, 85, medium  
**bad, 4, 81, light**  
bad, 6, 95, medium  
good, 4, 93, light

# Learning Binary Decision Trees

- Step 1: Start with all data points in a single node. Predict the most common value (classification) or mean value (regression).
- Step 2: Choose the “best” binary decision rule and use it to split the data into two groups.
- Step 3: Repeat step 2 on each group until some stopping criterion is reached.



MPG, cylinders, HP, weight

**good, 4, 75, light**

bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
**good, 4, 65, light**  
bad, 6, 85, medium  
**bad, 4, 81, light**  
bad, 6, 95, medium  
good, 4, 93, light

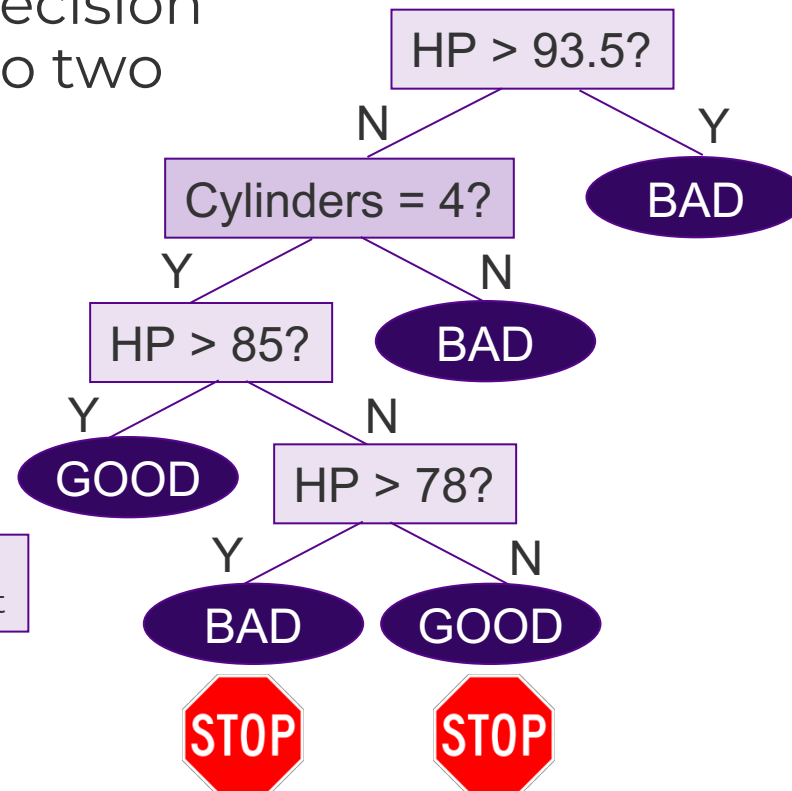
# Learning Binary Decision Trees

- Step 1: Start with all data points in a single node. Predict the most common value (classification) or mean value (regression).
- Step 2: Choose the “best” binary decision rule and use it to split the data into two groups.
- Step 3: Repeat step 2 on each group until some stopping criterion is reached.

- All outputs same?

- All inputs same?

bad, 4, 81, light  
good, 4, 81, light



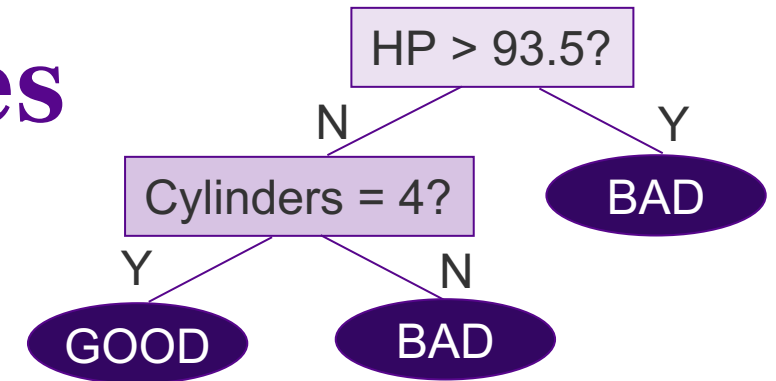
MPG, cylinders, HP, weight

**good, 4, 75, light**

bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
**good, 4, 65, light**  
bad, 6, 85, medium  
**bad, 4, 81, light**  
bad, 6, 95, medium  
good, 4, 93, light

# Learning Binary Decision Trees

- Step 1: Start with all data points in a single node. Predict the most common value (classification) or mean value (regression).
- Step 2: Choose the “best” binary decision rule and use it to split the data into two groups.
- Step 3: Repeat step 2 on each group until some stopping criterion is reached.
- Step 4: Prune the tree to remove irrelevant rules.



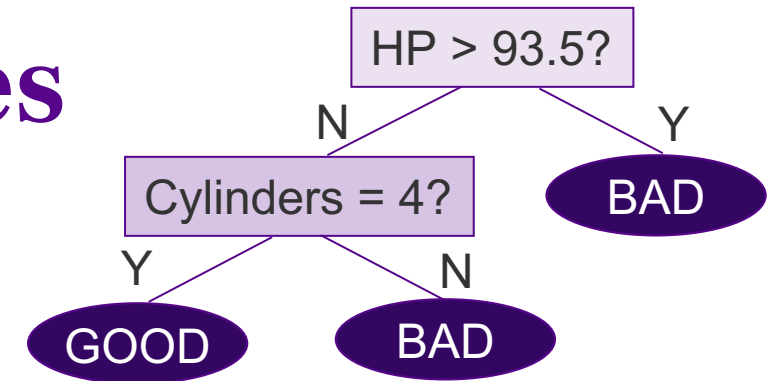
MPG, cylinders, HP, weight

good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
bad, 4, 81, light  
bad, 6, 95, medium  
good, 4, 93, light



# Learning Binary Decision Trees

- Step 1: Start with all data points in a single node. Predict the most common value (classification) or mean value (regression).
- Step 2: Choose the “best” binary decision rule and use it to split the data into two groups.
- Step 3: Repeat step 2 on each group until some stopping criterion is reached.
- Step 4: Prune the tree to remove irrelevant rules.



Question 1: How to choose the best decision rule for a given node?

Question 2: How to prune the tree (and why bother?)

MPG, cylinders, HP, weight

good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
bad, 4, 81, light  
bad, 6, 95, medium  
good, 4, 93, light

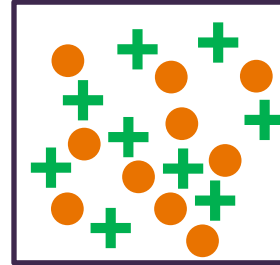
# Q1: How to Choose Best Split

# Q1: How to Choose Best Split

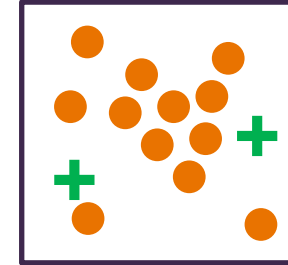
## Impurity (Diversity) Measures:

- Entropy (information gain)
- Information Gain Ratio
- Gini (population diversity)
- Sum of Square Errors

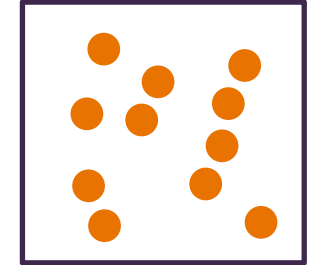
Very impure group



Less impure



Minimum impurity



Want to measure “purity” of the split

More certain about Yes/No after split:

- Pure set (**5 yes/15 no**), then completely certain (100%)
- Impure set (**10 yes/10 no**), then completely uncertain (50%)

Can't use Purity (“yes”| set):

- Must be symmetric: **5 yes/0 no** as pure as **0 yes/5 no**

# Q1: How to Choose Best Split

Impurity (Diversity) Measures:

- **Entropy (information gain)**
- Information Gain Ratio
- Gini (population diversity)
- Sum of Square Errors

# Choosing a Decision Rule

- We can use any input attribute to split.
  - If discrete: choose a class, split into = and ≠.
  - If real: choose a threshold, split into > and ≤.
- To choose a threshold for a real attribute, sort the values and use midpoints.

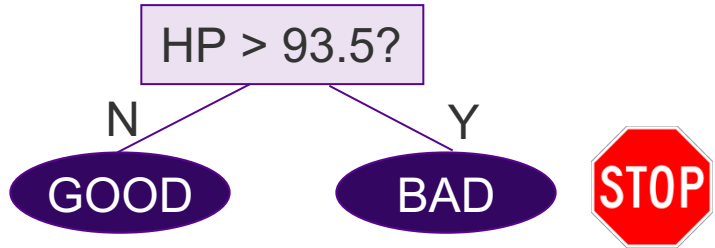
<u>Split</u>	<u>Group Y</u>	<u>Group N</u>
Cylinders = 4?	5+ / 4-	0+ / 11-
Cylinders = 6?	0+ / 6-	5+ / 9-
Cylinders = 8?	0+ / 5-	5+ / 10-
HP > 78?	2+ / 0-	3+ / 15-
HP > 87?	2+ / 2-	3+ / 13-
HP > 89.5?	3+ / 2-	2+ / 13-
HP > 91?	3+ / 3-	2+ / 12-
HP > 93.5?	5+ / 3-	0+ / 12-
Weight = light?	3+ / 3-	2+ / 12-
Weight = medium?	2+ / 6-	3+ / 9-
Weight = heavy?	0+ / 6-	5+ / 9-

MPG, cylinders, HP, weight

good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
bad, 4, 81, light  
bad, 6, 95, medium  
good, 4, 93, light

# Choosing a Decision Rule

- We can use any input attribute to split.
  - If discrete: choose a class, split into = and ≠.
  - If real: choose a threshold, split into > and ≤.
- To choose a threshold for a real attribute, sort the values and use midpoints.
- Choose the split with the highest information gain.



(5 good, 3 bad) (0 good, 12 bad)

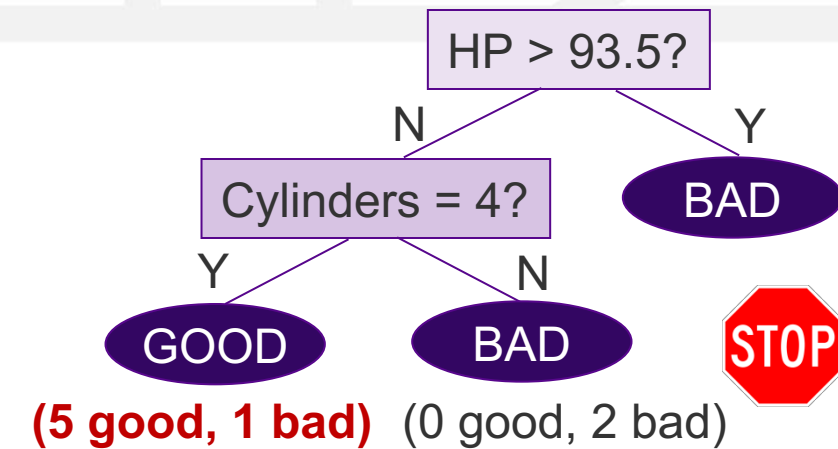
Split	Group Y	Group N	Gain
Cylinders = 4?	5+ / 4-	0+ / 11-	0.365
Cylinders = 6?	0+ / 6-	5+ / 9-	0.153
Cylinders = 8?	0+ / 5-	5+ / 10-	0.123
HP > 78?	2+ / 0-	3+ / 15-	0.226
HP > 87?	2+ / 2-	3+ / 13-	0.054
HP > 89.5?	3+ / 2-	2+ / 13-	0.144
HP > 91?	3+ / 3-	2+ / 12-	0.097
<b>HP &gt; 93.5?</b>	<b>5+ / 3-</b>	<b>0+ / 12-</b>	<b>0.430</b>
Weight = light?	3+ / 3-	2+ / 12-	0.097
Weight = medium?	2+ / 6-	3+ / 9-	0.000
Weight = heavy?	0+ / 6-	5+ / 9-	0.153

MPG, cylinders, HP, weight

good, 4, 75, light  
 bad, 6, 90, medium  
 bad, 4, 110, medium  
 bad, 8, 175, weighty  
 bad, 6, 95, medium  
 bad, 4, 94, light  
 bad, 4, 95, light  
 bad, 8, 139, weighty  
 bad, 8, 190, weighty  
 bad, 8, 145, weighty  
 bad, 6, 100, medium  
 good, 4, 92, medium  
 bad, 6, 100, weighty  
 bad, 8, 170, weighty  
 good, 4, 89, medium  
 good, 4, 65, light  
 bad, 6, 85, medium  
 bad, 4, 81, light  
 bad, 6, 95, medium  
 good, 4, 93, light

# Choosing a Decision Rule

- We repeat this process for each non-terminal node of the tree.



Split	Group Y	Group N	Gain
<b>Cylinders = 4?</b>	<b>5+ / 1-</b>	<b>0+ / 2-</b>	<b>0.467</b>
HP > 78?	2+ / 0-	3+ / 3-	0.204
HP > 87?	2+ / 2-	3+ / 1-	0.049
HP > 89.5?	3+ / 2-	2+ / 1-	0.003
HP > 91?	3+ / 3-	2+ / 0-	0.204
Weight = light?	3+ / 1-	2+ / 2-	0.049

MPG, cylinders, HP, weight

**good, 4, 75, light**

**bad, 6, 90, medium**

bad, 4, 110, medium

bad, 8, 175, weighty

bad, 6, 95, medium

bad, 4, 94, light

bad, 4, 95, light

bad, 8, 139, weighty

bad, 8, 190, weighty

bad, 8, 145, weighty

bad, 6, 100, medium

**good, 4, 92, medium**

bad, 6, 100, weighty

bad, 8, 170, weighty

**good, 4, 89, medium**

**good, 4, 65, light**

**bad, 6, 85, medium**

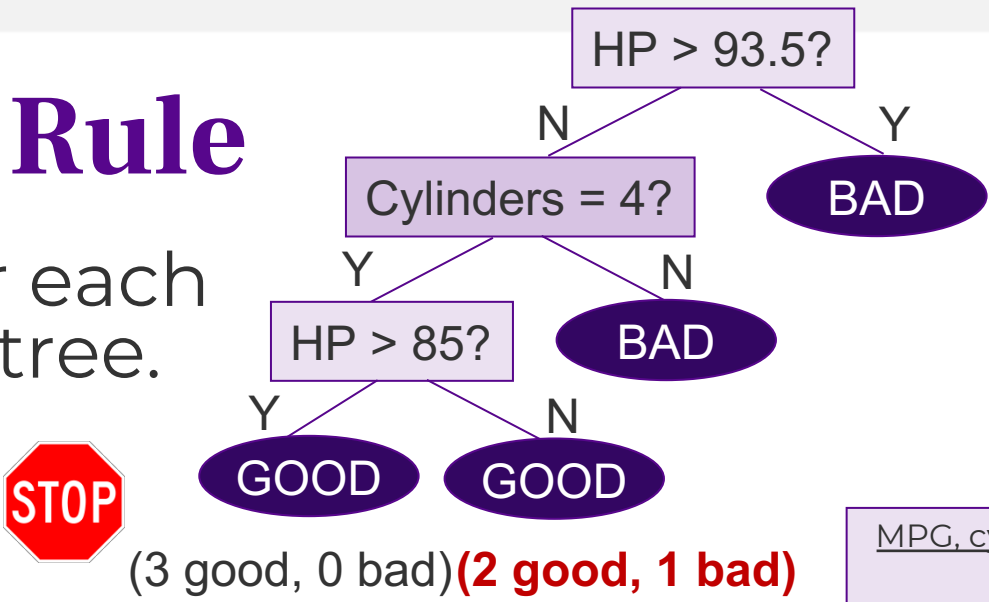
**bad, 4, 81, light**

bad, 6, 95, medium

**good, 4, 93, light**

# Choosing a Decision Rule

- We repeat this process for each non-terminal node of the tree.



Split	Group Y	Group N	Gain
HP > 78?	3+ / 1-	2+ / 0-	0.109
<b>HP &gt; 85?</b>	<b>3+ / 0-</b>	<b>2+ / 1-</b>	<b>0.191</b>
Weight = light?	3+ / 1-	2+ / 0-	0.109

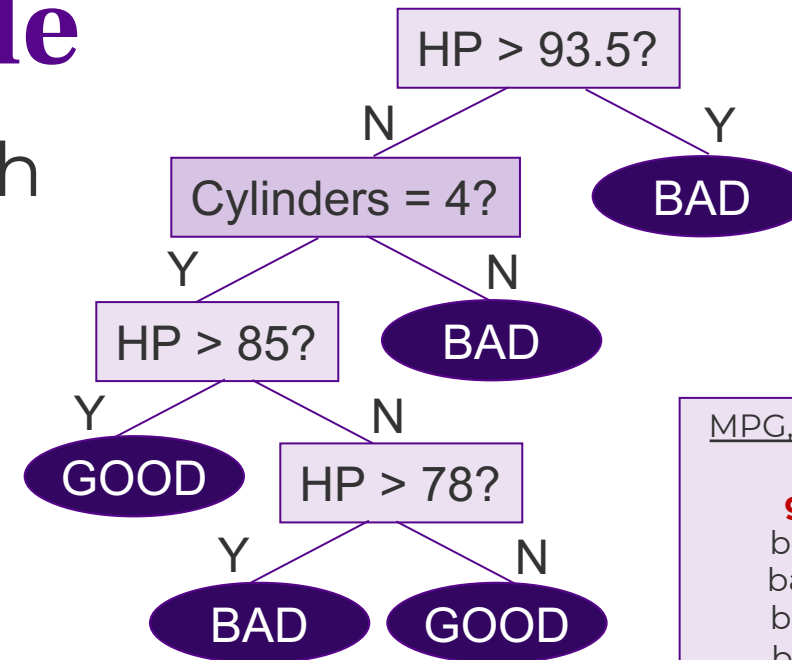
MPC, cylinders, HP, weight

**good, 4, 75, light**  
 bad, 6, 90, medium  
 bad, 4, 110, medium  
 bad, 8, 175, weighty  
 bad, 6, 95, medium  
 bad, 4, 94, light  
 bad, 4, 95, light  
 bad, 8, 139, weighty  
 bad, 8, 190, weighty  
 bad, 8, 145, weighty  
 bad, 6, 100, medium  
**good, 4, 92, medium**  
 bad, 6, 100, weighty  
 bad, 8, 170, weighty  
**good, 4, 89, medium**  
**good, 4, 65, light**  
 bad, 6, 85, medium  
**bad, 4, 81, light**  
 bad, 6, 95, medium  
**good, 4, 93, light**



# Choosing a Decision Rule

- We repeat this process for each non-terminal node of the tree.



<u>Split</u>	<u>Group Y</u>	<u>Group N</u>	<u>Gain</u>
<b>HP &gt; 78?</b>	<b>0+ / 1-</b>	<b>2+ / 0-</b>	<b>0.918</b>

MPG, cylinders, HP, weight

**good, 4, 75, light**

bad, 6, 90, medium

bad, 4, 110, medium

bad, 8, 175, weighty

bad, 6, 95, medium

bad, 4, 94, light

bad, 4, 95, light

bad, 8, 139, weighty

bad, 8, 190, weighty

bad, 8, 145, weighty

bad, 6, 100, medium

good, 4, 92, medium

bad, 6, 100, weighty

bad, 8, 170, weighty

good, 4, 89, medium

**good, 4, 65, light**

bad, 6, 85, medium

**bad, 4, 81, light**

bad, 6, 95, medium

good, 4, 93, light

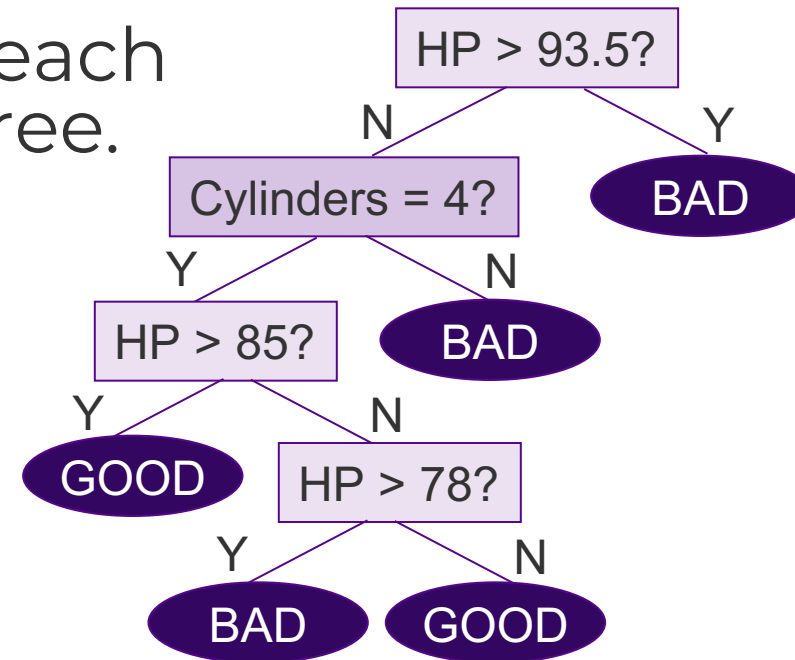
# Choosing a Decision Rule

- We repeat this process for each non-terminal node of the tree.

Information gain is an information-theoretic measure of how well the split separates the data.

It can be computed as a function of the numbers of + and – examples in each group.

<u>Group Y</u>	<u>Group N</u>
A+ / B-	C+ / D-



MPC, cylinders, HP, weight

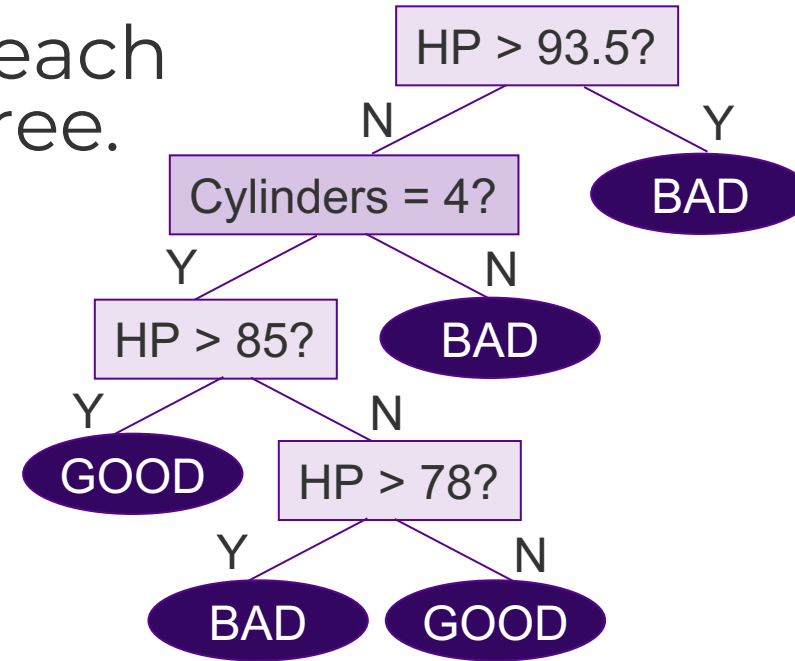
good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
bad, 4, 81, light  
bad, 6, 95, medium  
good, 4, 93, light

# Choosing a Decision Rule

- We repeat this process for each non-terminal node of the tree.

Information gain is an information-theoretic measure of how well the split separates the data.

It can be computed as a function of the numbers of + and – examples in each group.



Group Y	Group N
A+ / B-	C+ / D-

$$\text{Gain} = \frac{F((A + C), (B + D)) - F(A, B) - F(C, D)}{A + B + C + D}$$

$$\text{where: } F(X, Y) = X \log_2 \frac{X + Y}{X} + Y \log_2 \frac{X + Y}{Y}$$

*(You don't have to memorize this formula.)*

MPC, cylinders, HP, weight

good, 4, 75, light  
 bad, 6, 90, medium  
 bad, 4, 110, medium  
 bad, 8, 175, weighty  
 bad, 6, 95, medium  
 bad, 4, 94, light  
 bad, 4, 95, light  
 bad, 8, 139, weighty  
 bad, 8, 190, weighty  
 bad, 8, 145, weighty  
 bad, 6, 100, medium  
 good, 4, 92, medium  
 bad, 6, 100, weighty  
 bad, 8, 170, weighty  
 good, 4, 89, medium  
 good, 4, 65, light  
 bad, 6, 85, medium  
 bad, 4, 81, light  
 bad, 6, 95, medium  
 good, 4, 93, light

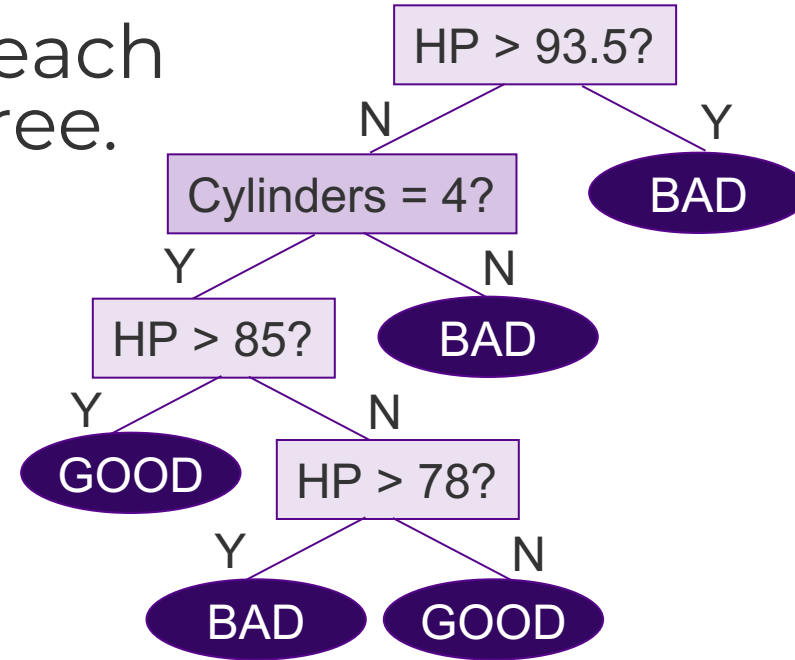
# Choosing a Decision Rule

- We repeat this process for each non-terminal node of the tree.

Information gain is an information-theoretic measure of how well the split separates the data.

It can be computed as a function of the numbers of + and – examples in each group.

<u>Group Y</u>	<u>Group N</u>	<u>Gain</u>
3+ / 1-	6+ / 2-	0.000
8+ / 1-	1+ / 2-	0.204
9+ / 0-	0+ / 3-	0.811



Intuitively, the information gain is large when the proportions of positive examples in the two groups are very different and zero when they are the same.

MPC, cylinders, HP, weight

good, 4, 75, light  
 bad, 6, 90, medium  
 bad, 4, 110, medium  
 bad, 8, 175, weighty  
 bad, 6, 95, medium  
 bad, 4, 94, light  
 bad, 4, 95, light  
 bad, 8, 139, weighty  
 bad, 8, 190, weighty  
 bad, 8, 145, weighty  
 bad, 6, 100, medium  
 good, 4, 92, medium  
 bad, 6, 100, weighty  
 bad, 8, 170, weighty  
 good, 4, 89, medium  
 good, 4, 65, light  
 bad, 6, 85, medium  
 bad, 4, 81, light  
 bad, 6, 95, medium  
 good, 4, 93, light

# Q1: How to Choose Best Split

## Impurity (Diversity) Measures:

- Entropy (information gain)
- **Information Gain Ratio**
- Gini (population diversity)
- Sum of Square Errors

The solution to this problem is used to somehow penalize the attributes that lead to a very high number of branches.

- One option is to take into account the number and size of the children nodes, regardless of what classes they contain.
- Instead of using “gain” to determine which attribute to use for the next branch, we shall use “gain ratio,” which we shall define as the division between gain and the information value of the attribute.

$$\text{GainRatio}(T) = \text{Gain}(T) / \text{SplitInfo}(T)$$

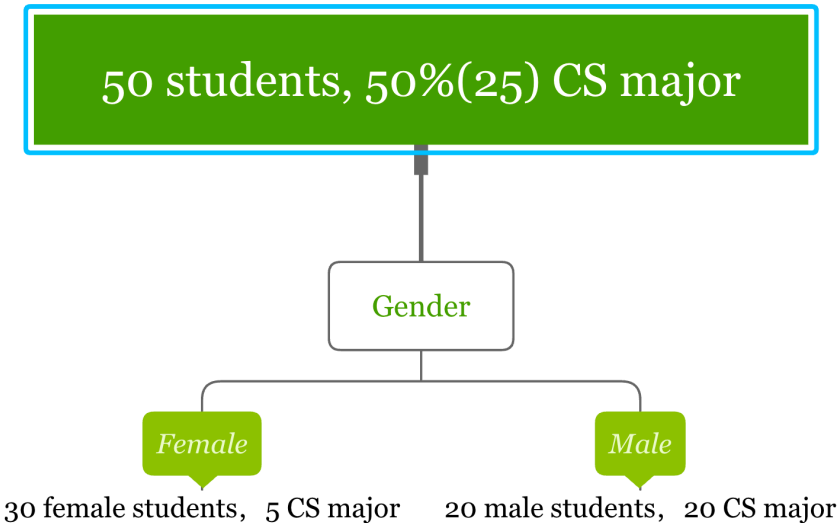
# Q1: How to Choose Best Split

Gini impurity is a measure of non-homogeneity.

$$\text{Gini Index} = \sum_i p_i(1 - p_i)$$

## Impurity (Diversity) Measures:

- Entropy (information gain)
- Information Gain Ratio
- **Gini (population diversity)**
- Sum of Square Errors



Let's calculate the Gini impurity for splitting node "Gender":

1. Gini impurity for "Female" =  $16 \times 56 + 56 \times 16 = 518$
2. Gini impurity for "Male" =  $0 \times 1 + 1 \times 0 = 0$

The Gini impurity for the node "Gender" is the following weighted average of the above two scores:  $35 \times 518 + 25 \times 0 = 16$

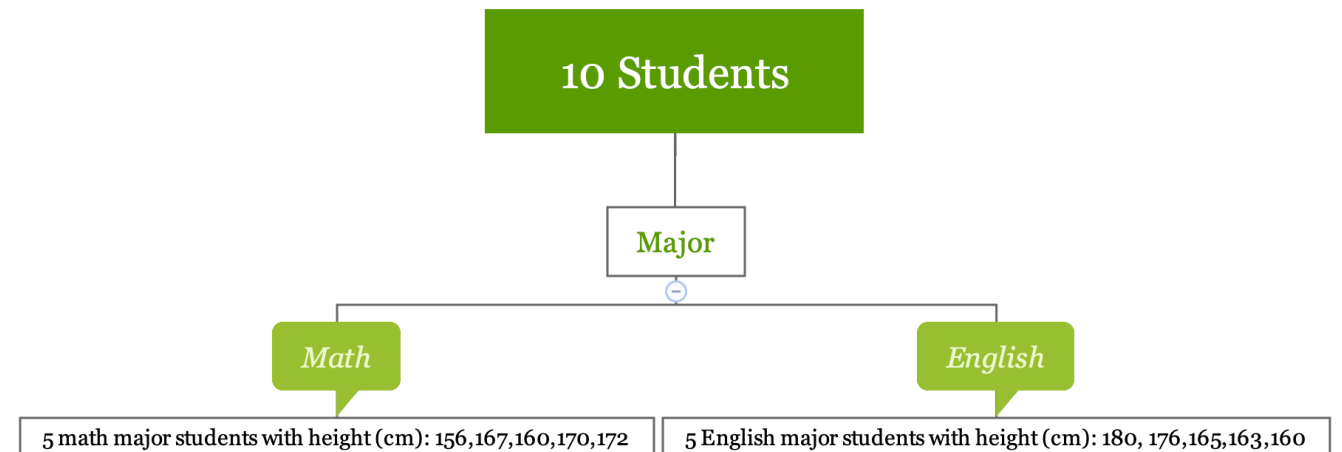
# Q1: How to Choose Best Split

The SSE is the most widely used splitting metric for regression.

## Impurity (Diversity) Measures:

- Entropy (information gain)
- Information Gain Ratio
- Gini (population diversity)
- **Sum of Square Errors**

$$SSE = \sum_{i \in S_1} (y_i - \bar{y}_1)^2 + \sum_{i \in S_2} (y_i - \bar{y}_2)^2$$



In this situation:

1. SSE for “Math” is 184
2. SSE for “English” is 302.8
3. SSE for splitting node “Major” is the sum of the above two numbers, which is 486.8

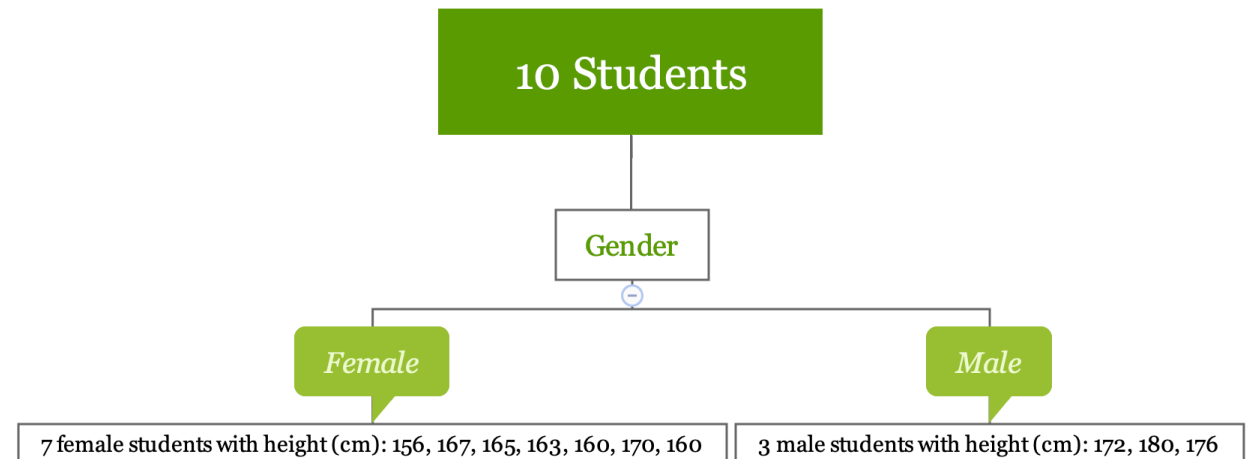
# Q1: How to Choose Best Split

The SSE is the most widely used splitting metric for regression.

## Impurity (Diversity) Measures:

- Entropy (information gain)
- Information Gain Ratio
- Gini (population diversity)
- **Sum of Square Errors**

$$SSE = \sum_{i \in S_1} (y_i - \bar{y}_1)^2 + \sum_{i \in S_2} (y_i - \bar{y}_2)^2$$



In this situation:

1. SSE for “Female” is 136
2. SSE for “Male” is 32
3. SSE for splitting node “Gender” is the sum of the above two numbers which is 168
4. SSE for the 10 students in root node is 522.9. After the split, SSE decreases from 522.9 to 168.

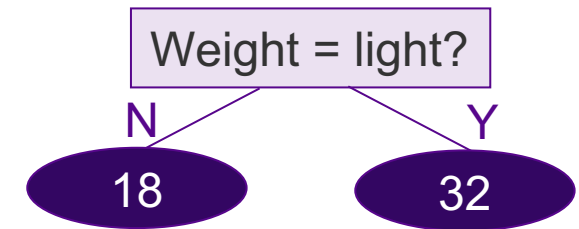


# Choosing a Decision Rule

- If we are trying to predict a real-valued attribute (i.e. doing regression), minimize the sum of squared errors instead of maximizing information gain.

MPG, cylinders, HP, weight

32, 4, 75, light  
20, 6, 95, medium  
20, 4, 115, medium  
14, 6, 95, medium



<u>Split</u>	<u>Group Y</u>	<u>Group N</u>	<u>Total SSE</u>
Cylinders = 4?	Predict: 26 SSE: 72	Predict: 17 SSE: 18	90
<b>Weight = light?</b>	<b>Predict: 32</b> <b>SSE: 0</b>	<b>Predict: 18</b> <b>SSE: 24</b>	<b>24</b>
<b>HP &gt; 85?</b>	<b>Predict: 18</b> <b>SSE: 24</b>	<b>Predict: 32</b> <b>SSE: 0</b>	<b>24</b>
HP > 105?	Predict: 20 SSE: 0	Predict: 22 SSE: 168	168

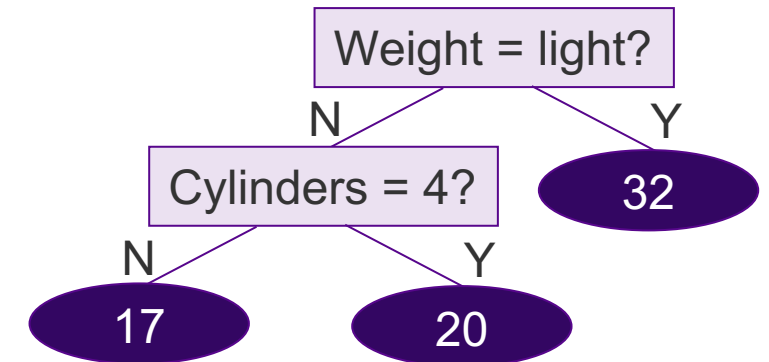
# Choosing a Decision Rule

- If we are trying to predict a real-valued attribute (i.e. doing regression), minimize the sum of squared errors instead of maximizing information gain.

<u>Split</u>	<u>Group Y</u>	<u>Group N</u>	<u>Total SSE</u>
<b>Cylinders = 4?</b>	<b>Predict: 20</b> <b>SSE: 0</b>	<b>Predict: 17</b> <b>SSE: 18</b>	<b>18</b>
<b>HP &gt; 105?</b>	<b>Predict: 20</b> <b>SSE: 0</b>	<b>Predict: 17</b> <b>SSE: 18</b>	<b>18</b>

MPG, cylinders, HP, weight

32, 4, 75, light  
20, 6, 95, medium  
20, 4, 115, medium  
14, 6, 95, medium



# Mathematical Formulation

1

Given training vectors  $x_i$  and a label vector  $y$ . Let data at tree node  $m$  be represented by  $Q$ .

$$x_i \in R^n \quad y \in R^l$$

2

Consider a set of candidate binary splits, each consisting of a feature  $j$  and threshold  $t_m$ , and partitioning  $Q$  into subsets  $Q_{\text{left}}$  and  $Q_{\text{right}}$ .

$$\theta = (j, t_m)$$

$$Q_{\text{left}}(\theta) = (x, y) | x_j \leq t_m$$

$$Q_{\text{right}}(\theta) = Q \setminus Q_{\text{left}}(\theta)$$

3

Select the split that minimizes the average *impurity* of  $Q_{\text{left}}$  and  $Q_{\text{right}}$ , and recurse.

$$\theta^* = \operatorname{argmin}_{\theta} G(Q, \theta)$$

$$G(Q, \theta) = \frac{n_{\text{left}}}{N_m} H(Q_{\text{left}}(\theta)) + \frac{n_{\text{right}}}{N_m} H(Q_{\text{right}}(\theta))$$

4<sup>a</sup>

For classification, use *cross-entropy*. ( $p_{mk}$  are class proportions at node  $m$ ).

$$p_{mk} = 1/N_m \sum_{x_i \in R_m} I(y_i = k)$$

$$H(X_m) = - \sum_k p_{mk} \log(p_{mk})$$

4<sup>b</sup>

For regression, use *mean squared error* ( $c_m$  is the mean at node  $m$ ).

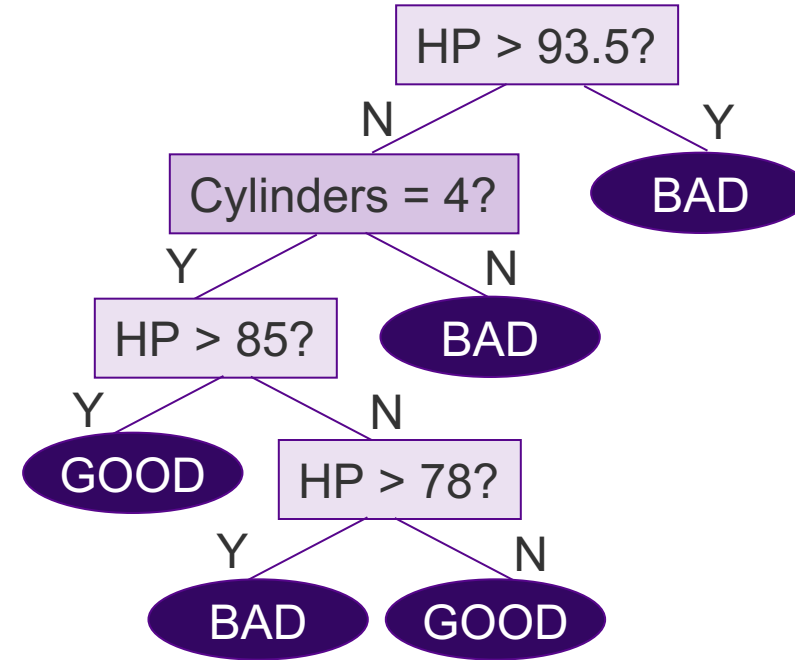
$$c_m = \frac{1}{N_m} \sum_{i \in N_m} y_i$$

$$H(X_m) = \frac{1}{N_m} \sum_{i \in N_m} (y_i - c_m)^2$$

## Q2: How to Prune the Tree

# Pruning Decision Trees

- Notice that the unpruned decision tree classifies every training example perfectly.
- This will always be the case (unless there are records with the same input values and different output values, as in the regression example).
- Does this mean we've found the best tree?



MPG, cylinders, HP, weight

good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
bad, 4, 81, light  
bad, 6, 95, medium  
good, 4, 93, light

What we really want to know:

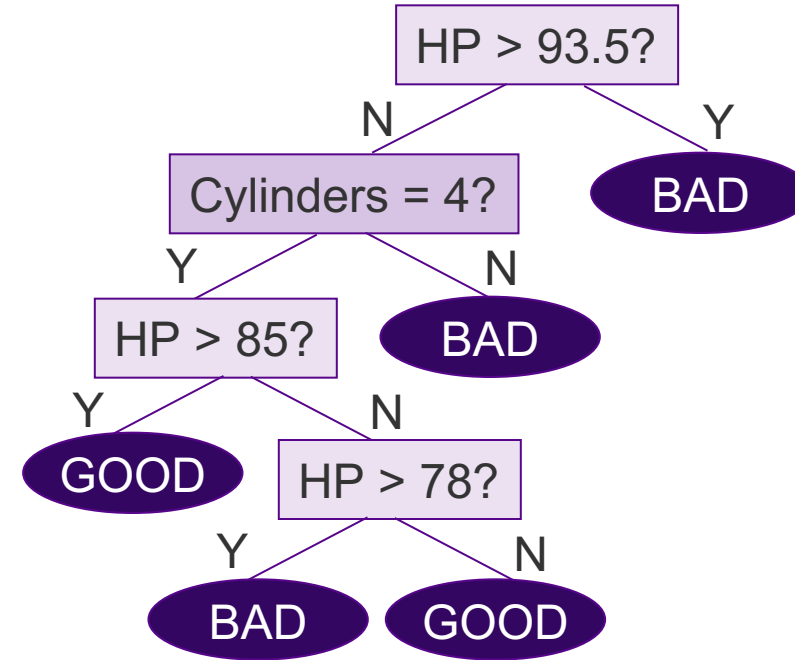
How well does this classifier predict values for new data that we haven't already seen?

# Pruning Decision Trees

- One way to answer this question:
  - Hide part of your data (the “test set”).
  - Learn a tree using the rest of the data (the “training set”).
  - See what proportion of the test set is classified correctly.
- Consider pruning each node to see if it reduces test set error.

Training set (20 examples):  
**100%** correct, **0%** incorrect

Test set (100 examples):  
**86%** correct, **14%** incorrect

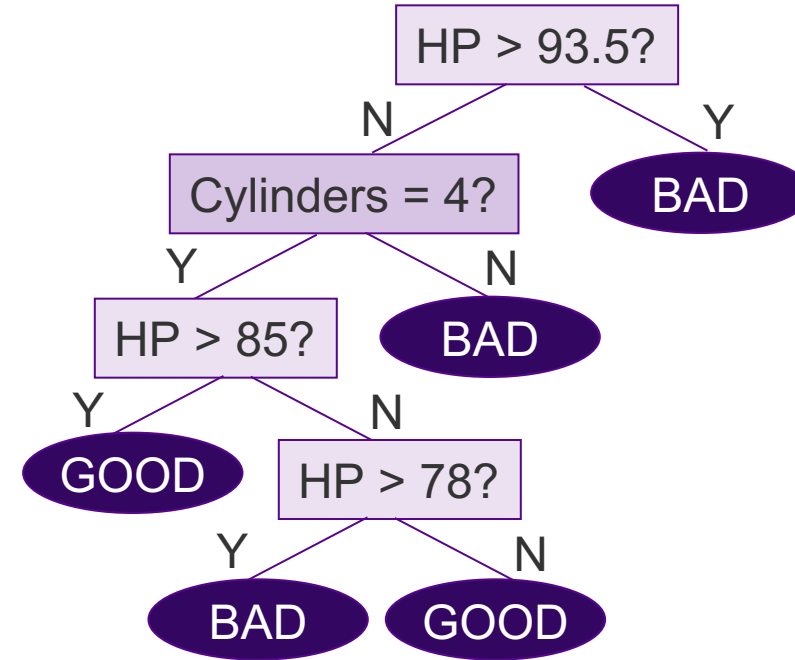


MPG, cylinders, HP, weight

good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
bad, 4, 81, light  
bad, 6, 95, medium  
good, 4, 93, light

# Pruning Decision Trees

- One way to answer this question:
  - Hide part of your data (the “test set”).
  - Learn a tree using the rest of the data (the “training set”).
  - See what proportion of the test set is classified correctly.
- Consider pruning each node to see if it reduces test set error.



Training set (20 examples):  
**100%** correct, **0%** incorrect

Test set (100 examples):  
**86%** correct, **14%** incorrect

MPG, cylinders, HP, weight

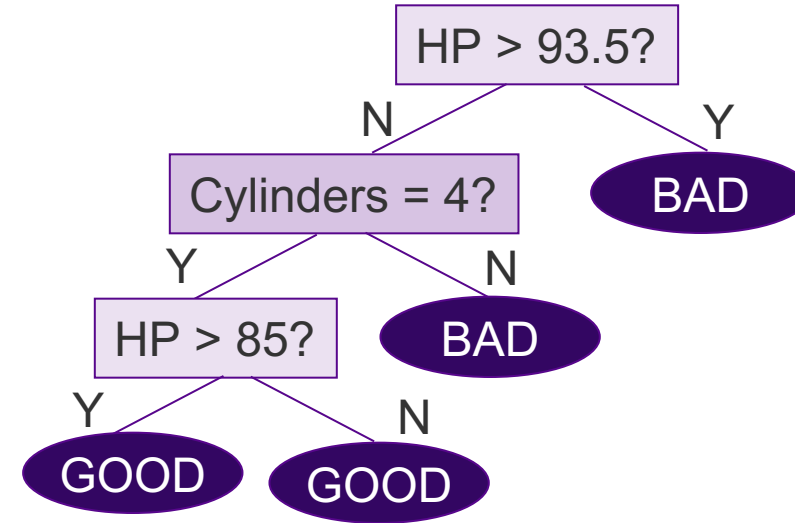
good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
bad, 4, 81, light  
bad, 6, 95, medium  
good, 4, 93, light

What we really want to know:

How well does this classifier predict values for new data that we haven't already seen?

# Pruning Decision Trees

- One way to answer this question:
  - Hide part of your data (the “test set”).
  - Learn a tree using the rest of the data (the “training set”).
  - See what proportion of the test set is classified correctly.
- Consider pruning each node to see if it reduces test set error.



Training set (20 examples):  
**95%** correct, **5%** incorrect

Test set (100 examples):  
**89%** correct, **11%** incorrect

MPG, cylinders, HP, weight

good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
**bad, 4, 81, light**  
bad, 6, 95, medium  
good, 4, 93, light

What we really want to know:  
How well does this classifier predict values for new data  
that we haven't already seen?

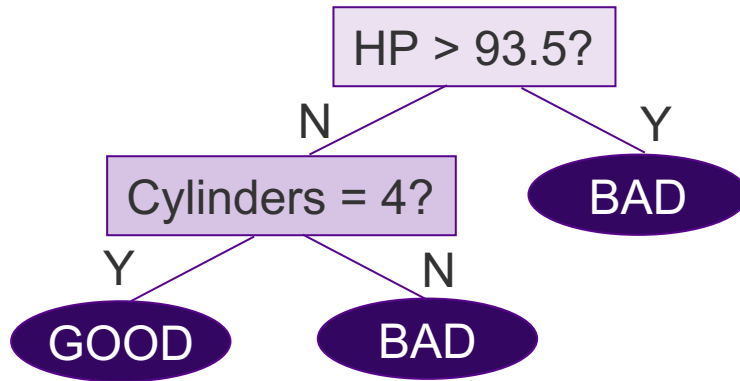


# Pruning Decision Trees

- One way to answer this question:
  - Hide part of your data (the “test set”).
  - Learn a tree using the rest of the data (the “training set”).
  - See what proportion of the test set is classified correctly.
- Consider pruning each node to see if it reduces test set error.

Training set (20 examples):  
**95%** correct, **5%** incorrect

Test set (100 examples):  
**89%** correct, **11%** incorrect



MPG, cylinders, HP, weight

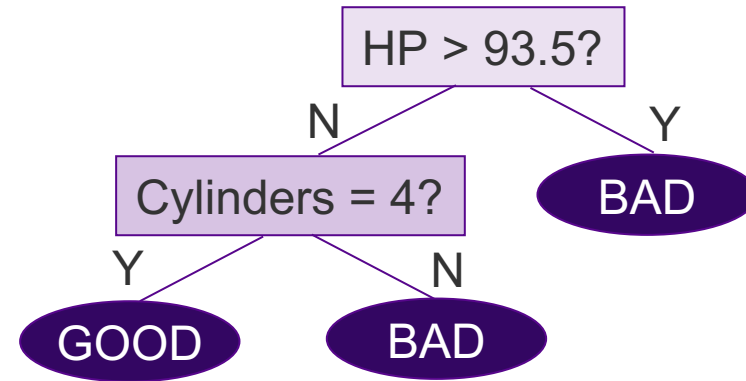
good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
**bad, 4, 81, light**  
bad, 6, 95, medium  
good, 4, 93, light

What we really want to know:

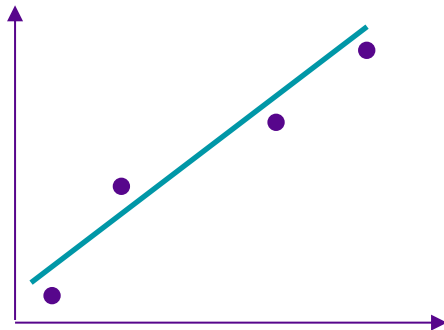
How well does this classifier predict values for new data that we haven't already seen?

# Pruning Decision Trees

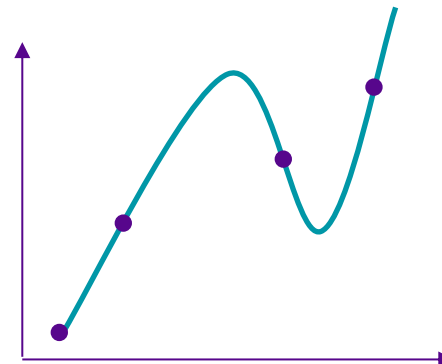
- Q: Why does pruning the tree reduce test set error?
- A: Because the unpruned tree was overfitting the training data (paying attention to parts of the data that are not relevant for prediction).



Here's another example of overfitting, this time for regression:



A line fits pretty well...



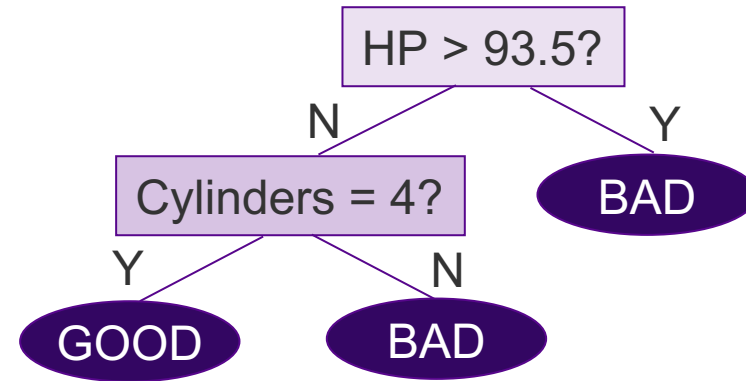
A cubic is probably overfitting.

MPG, cylinders, HP, weight

good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
**bad, 4, 81, light**  
bad, 6, 95, medium  
good, 4, 93, light

# Pruning Decision Trees

- Q: What if we don't have enough data points?
- A: Another way to prevent overfitting is to do a certain kind of significance test (chi-squared) for each node and prune any nodes that are not significant.



MPG, cylinders, HP, weight

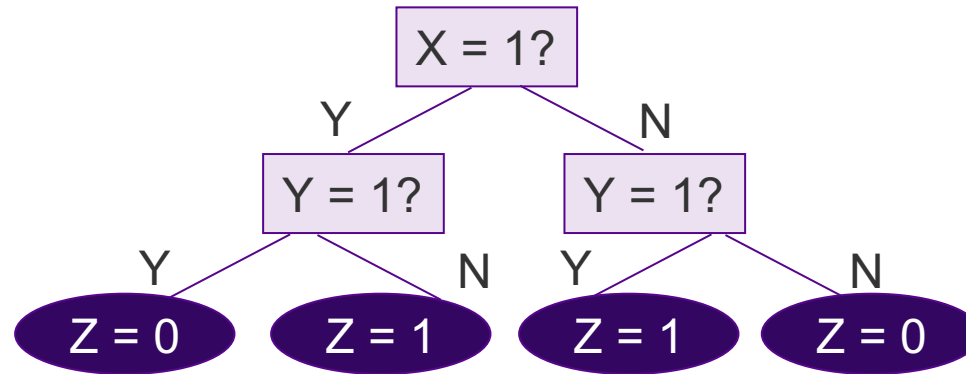
good, 4, 75, light  
bad, 6, 90, medium  
bad, 4, 110, medium  
bad, 8, 175, weighty  
bad, 6, 95, medium  
bad, 4, 94, light  
bad, 4, 95, light  
bad, 8, 139, weighty  
bad, 8, 190, weighty  
bad, 8, 145, weighty  
bad, 6, 100, medium  
good, 4, 92, medium  
bad, 6, 100, weighty  
bad, 8, 170, weighty  
good, 4, 89, medium  
good, 4, 65, light  
bad, 6, 85, medium  
**bad, 4, 81, light**  
bad, 6, 95, medium  
good, 4, 93, light

# Pruning Decision Trees

- Q: Why bother building the whole tree if we're just going to prune it?
- A: We could do significance testing while building the tree, but for many datasets, post-pruning gives better performance.

Consider the XOR function:

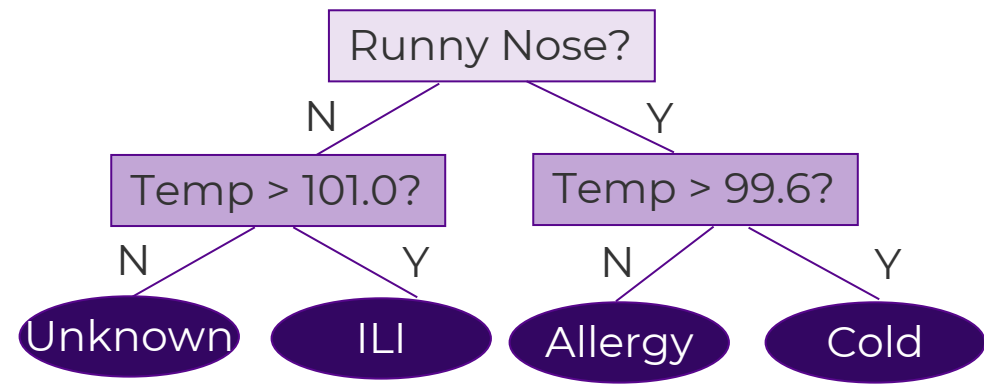
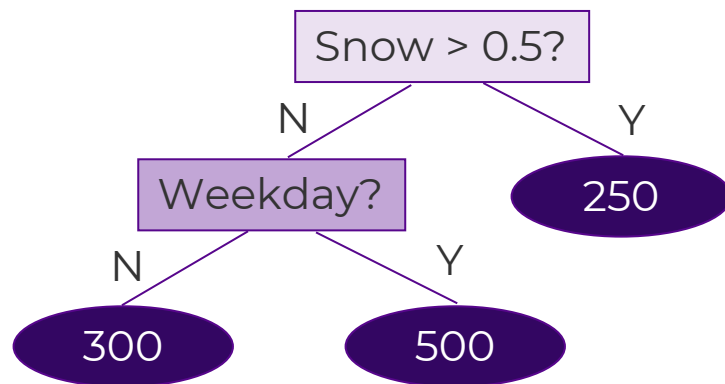
<u>X</u>	<u>Y</u>	<u>Z</u>
0	0	0
0	1	1
1	0	1
1	1	0



Each second-level split has high information gain. The first level split has no information gain but is needed to make the second level splits possible.

# Some Advantages of Decision Trees

- Easy to learn the tree automatically from a dataset.
- Very easy to predict the target value given the tree.
- Generally good classification performance (though some fancier methods may do better, depending on the dataset).
- It can do both classification and regression and can use both real and discrete inputs.
- It gives an idea of which variables are important in predicting the target value.
  - More important variables tend to be toward the top of the tree.
  - Unimportant variables are not included.



# Some Issues with Decision Trees

1) Multi-way instead of binary splits are okay, but we generally prefer binary.

If using multi-way, use gain ratio instead of information gain to adjust for the arity (# of distinct values) of an attribute. Otherwise, the dataset gets fragmented, leading to poor generalization accuracy (e.g., splitting on country name in the illness dataset).

2) Decision trees can be easily used to predict multiple outputs.

Create a single tree where the split criterion is the average impurity over all outputs.

3) Lots of different ways to handle missing data:

(a) Delete missing observations (not so good); (b) Treat “missing” as its own value (good only if missingness is informative); (c) Surrogate splits (Duda et al. Ch. 8.3.10); (d) Propagate examples with missing values down both branches as partial observations

4) Can create a rule list from the leaves of a decision tree.

Advantages: possibly more interpretable; can prune the rule list directly.

# Lab Time

# For the Next Week (Week 3)

1. Read Python documentation at  
<http://scikit-learn.org/stable/modules/tree.html>
2. Check readings and review them



# References

1. L. Breiman, J.H. Friedman, R.A. Olshen, and C.J. Stone. Classification and Regression Trees, Wadsworth, 1984.
2. J.R. Quinlan. C4.5 : Programs for Machine Learning, Morgan Kaufmann, 1993.
3. T.M. Mitchell. Machine Learning, McGraw-Hill, 1997. (Chapter 3)
4. T. Hastie, R. Tibshirani, J. Friedman. The Elements of Statistical Learning, 2001.