Decision Trees

Nipun Batra and teaching staff July 26, 2025

IIT Gandhinagar

Table of Contents

Introduction and Motivation

Discrete Input, Real Output

Real Input Discrete Output

Real Input Real Output

Pruning and Overfitting

Summary and Key Takeaways

Weighted Entropy

Introduction and Motivation

The need for interpretability

How to maintain trust in AI

Beyond developing initial trust, however, creators of AI also must work to maintain that trust. Siau and Wang suggest seven ways of "developing continuous trust" beyond the initial phases of product development:

- Usability and reliability. AI "should be designed to operate easily and intuitively," Siau and Wang write. "There should be no unexpected downtime or crashes."
- Collaboration and communication. AI developers want to create systems that perform autonomously, without human involvement. Developers must focus on creating AI applications that smoothly and easily collaborate and communicate with humans.
- Sociability and bonding. Building social activities into AI applications is one way to strengthen trust. A robotic dog that can recognize its owner and show affection is one example, Siau and Wang write.
- Security and privacy protection. Al applications rely on large data sets, so
 ensuring privacy and security will be crucial to establishing trust in the
 applications.
- Interpretability. Just as transparency is instrumental in building initial trust, interpretability – or the ability for a machine to explain its conclusions or actions – will help sustain trust.

Training Data

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Learning a Complicated Neural Network



Learnt Decision Tree



Medical Diagnosis using Decision Trees



Source: Improving medical decision trees by combining relevant health-care criteria

Leo Brieman

CONTRACTOR -	Leo Breiman 1928-2005		FOLLOW	Cited by		VIEW ALL
	Professor of Statistics, <u>UC Berkeley</u>				All	Since 2015
	Data Analysis Statistics Machine Learning			Citations h-index i10-index	142857 51 80	68736 33 46
TITLE		CITED BY	YEAR			17000
Random forests L Breiman Machine learning 45 (1), 5-32	53816	2001		пH	12750
Classification and L Breiman, JH Friedma CRC Press, New York	Regression Trees In, RA Olshen, CJ Stone	43992 *	1999	111		4250
Classification and L Breiman Chapman & Hall/CRC	regression trees	43992 *	1984	2013 2014 2015	2016 2017 2018 20	19 2020 0
Bagging predictors L Breiman Machine learning 24 (2	s 1), 123-140	22742	1996			
Statistical Modeling	g: The Two Cutures	2788 *	2003			
Statistical modeling L Breiman Statistical Science 16 (g: The two cultures (with comments and a rejoinder by the author) 3), 199-231	2772	2001			
Estimating optimal	transformations for multiple regression and correlation	2096	1985			

Optimal Decision Tree

Volume 5, number 1

INFORMATION PROCESSING LETTERS

May 1976

CONSTRUCTING OPTIMAL BINARY DECISION TREES IS NP-COMPLETE*

Laurent HYAFIL IRIA – Laboria, 78150 Rocquencourt, France

and

Ronald L. RIVEST Dept. of Electrical Engineering and Computer Science, M.I.T., Cambridge, Massachusetts 02139, USA

Received 7 November 1975, revised version received 26 January 1976

Binary decision trees, computational complexity, NP-complete

Quick Question!

Why is finding the optimal decision tree NP-hard?

A) The number of possible trees grows exponentially with features

Quick Question!

- A) The number of possible trees grows exponentially with features
- B) We need to consider all possible splits at each node

Quick Question!

- A) The number of possible trees grows exponentially with features
- B) We need to consider all possible splits at each node
- C) The problem requires checking all subsets of training data

Quick Question!

- A) The number of possible trees grows exponentially with features
- B) We need to consider all possible splits at each node
- C) The problem requires checking all subsets of training data
- D) All of the above

Quick Question!

- A) The number of possible trees grows exponentially with features
- B) We need to consider all possible splits at each node
- C) The problem requires checking all subsets of training data
- D) All of the above

Quick Question!

Why is finding the optimal decision tree NP-hard?

- A) The number of possible trees grows exponentially with features
- B) We need to consider all possible splits at each node
- C) The problem requires checking all subsets of training data
- D) All of the above

Answer: D) All of the above - The search space is exponentially large, making brute force optimization computationally intractable.

Core idea: At each level, choose an attribute that gives **biggest** estimated performance gain!



Image source: analyticsvidhya

 $\mathsf{Greedy} \neq \mathsf{Optimal}$

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

• For examples, we have 9 Yes, 5 No

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

• For examples, we have 9 Yes, 5 No

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?

Yes!

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?

Yes!

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?
- Yes!
- Key insight: Problem is "easier" when there is less disagreement

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?
- Yes!
- Key insight: Problem is "easier" when there is less disagreement

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- For examples, we have 9 Yes, 5 No
- Would it be trivial if we had 14 Yes or 14 No?
- Yes!
- Key insight: Problem is "easier" when there is less disagreement
- Need some statistical measure of "disagreement"

Entropy

Statistical measure to characterize the (im)purity of examples

Entropy

Statistical measure to characterize the (im)purity of examples $H(X) = -\sum_{i=1}^{k} p(x_i) \log_2 p(x_i)$



Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No
Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

• Can we use Outlook as the root node?

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

• Can we use Outlook as the root node?

Towards biggest estimated performance gain

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- Can we use Outlook as the root node?
- When Outlook is overcast, we always Play and thus no "disagreement"

Reduction in entropy by partitioning examples (S) on attribute A

$$\mathsf{Gain}(S, A) \equiv \mathsf{Entropy}(S) - \sum_{v \in \mathsf{Values}(A)} \frac{|S_v|}{|S|} \mathsf{Entropy}(S_v)$$

Pop Quiz #7

Quick Question!

What does entropy measure in the context of decision trees?

A) The depth of the tree

- A) The depth of the tree
- B) The impurity or "disagreement" in a set of examples

- A) The depth of the tree
- B) The impurity or "disagreement" in a set of examples
- C) The number of features in the dataset

- A) The depth of the tree
- B) The impurity or "disagreement" in a set of examples
- C) The number of features in the dataset
- D) The accuracy of the tree

Pop Quiz #11

Quick Question!

- A) The depth of the tree
- B) The impurity or "disagreement" in a set of examples
- C) The number of features in the dataset
- D) The accuracy of the tree

What does entropy measure in the context of decision trees?

- A) The depth of the tree
- B) The impurity or "disagreement" in a set of examples
- C) The number of features in the dataset
- D) The accuracy of the tree

Answer: B) The impurity or "disagreement" in a set of examples - Higher entropy means more mixed classes, lower entropy means more pure subsets.

• Create a root node for tree

- Create a root node for tree
- If all examples are +/-, return root with label = +/-

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - $\bullet~\mathsf{A} \leftarrow \mathsf{attribute}$ from Attributes which best classifies Examples

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - A \leftarrow attribute from Attributes which best classifies Examples
 - $\bullet \ \mathsf{Root} \leftarrow \mathsf{A}$

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - A \leftarrow attribute from Attributes which best classifies Examples
 - Root $\leftarrow A$
 - For each value (v) of A

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - A \leftarrow attribute from Attributes which best classifies Examples
 - Root $\leftarrow A$
 - For each value (v) of A
 - Add new tree branch : $\mathsf{A}=\mathsf{v}$

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - A \leftarrow attribute from Attributes which best classifies Examples
 - Root $\leftarrow A$
 - For each value (v) of A
 - Add new tree branch : A = v
 - Examples_v: subset of examples that A = v

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - A \leftarrow attribute from Attributes which best classifies Examples
 - Root $\leftarrow A$
 - For each value (v) of A
 - Add new tree branch : A = v
 - $\bullet \ \mathsf{Examples}_v : \ \mathsf{subset} \ \mathsf{of} \ \mathsf{examples} \ \mathsf{that} \ \mathsf{A} = \mathsf{v}$
 - If $Examples_v$ is empty: add leaf with label = most common value of Target Attribute

- Create a root node for tree
- If all examples are +/-, return root with label = +/-
- If attributes = empty, return root with most common value of Target Attribute in Examples
- Begin
 - A \leftarrow attribute from Attributes which best classifies Examples
 - Root $\leftarrow A$
 - For each value (v) of A
 - Add new tree branch : A = v
 - $\bullet \ \mathsf{Examples}_v : \ \mathsf{subset} \ \mathsf{of} \ \mathsf{examples} \ \mathsf{that} \ \mathsf{A} = \mathsf{v}$
 - If Examples_vis empty: add leaf with label = most common value of Target Attribute
 - Else: ID3 (Examples, Target attribute, Attributes A)

Learnt Decision Tree

Root Node (empty)

Training Data

Day	Outlook	Temp	Humidity	Windy	Play
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

We have 14 examples in S: 5 No, 9 Yes

Entropy(S) =
$$-p_{No} \log_2 p_{No} - p_{Yes} \log_2 p_{Yes}$$

= $-\frac{5}{14} \log_2 \left(\frac{5}{14}\right) - \frac{9}{14} \log_2 \left(\frac{9}{14}\right) = 0.940$

Outlook	Play
Sunny	No
Sunny	No
Overcast	Yes
Rain	Yes
Rain	Yes
Rain	No
Overcast	Yes
Sunny	No
Sunny	Yes
Rain	Yes
Sunny	Yes
Overcast	Yes
Overcast	Yes
Rain	No

Outlook	Play
Sunny	No
Sunny	No
Sunny	No
Sunny	Yes
Sunny	Yes
We have 2 \	es, 3 No

 $\begin{array}{l} {\sf Entropy} = \\ -\frac{3}{5}\log_2\left(\frac{3}{5}\right) - \\ \frac{2}{5}\log_2\left(\frac{2}{5}\right) = 0.971 \end{array}$

Outlook	Play
Sunny	No
Sunny	No
Sunny	No
Sunny	Yes
Sunny	Yes
We have 2 \	es, 3 No

 $\begin{array}{l} {\sf Entropy} = \\ -\frac{3}{5}\log_2\left(\frac{3}{5}\right) - \\ \frac{2}{5}\log_2\left(\frac{2}{5}\right) = 0.971 \end{array}$

Outlook	Play				
Sunny	No		Outlook	Play	
Sunny	No		Overcast	Yes	
Sunny	No		Overcast	Yes	
Sunny	Yes		Overcast	Yes	
Sunny	Yes		Overcast	Yes	
We have 2 Y	′es, 3 No	V	Ve have 4 Y	és, 0 No	
Entropy =			Entropy = 0 (pure		
$-\frac{3}{5}\log_2\left(\frac{3}{5}\right) -$			subset)		
$\frac{2}{5}\log_2\left(\frac{2}{5}\right) =$	= 0.971			-	

Outlook	Play				
Sunny	No		Outlook	Play	
Sunny	No		Overcast	Yes	
Sunny	No		Overcast	Yes	
Sunny	Yes		Overcast	Yes	
Sunny	Yes		Overcast	Yes	
We have 2 Y	′es, 3 No	V	Ve have 4 Y	és, 0 No	
Entropy =			Entropy = 0 (pure		
$-\frac{3}{5}\log_2\left(\frac{3}{5}\right) -$			subset)		
$\frac{2}{5}\log_2\left(\frac{2}{5}\right) =$	= 0.971			-	

Outlook	Play					Outlook	Play
Sunny	No		Outlook	Play		Rain	Yes
Sunny	No		Overcast	Yes		Rain	Yes
Sunny	No		Overcast	Yes		Rain	No
Sunny	Yes		Overcast	Yes		Rain	Yes
Sunny	Yes		Overcast	Yes		Rain	No
We have 2 Yes, 3 No		o \	We have 4 Yes, 0 No		٧	Ve have 3 ۲	′es, 2 No
Entropy =			Entropy = 0 (pure			Entrop	y =
$-\frac{3}{5}\log_2\left(\frac{3}{5}\right) -$			subset)			$-\frac{3}{5}\log_{2}($	$(\frac{3}{5}) -$
$\frac{2}{5}\log_2\left(\frac{2}{5}\right) =$	= 0.971				$\frac{2}{5}\log_2\left(\frac{2}{5}\right) = 0.971$		

$$\mathsf{Gain}(S,\mathsf{Outlook}) = \mathsf{Entropy}(S) - \sum_{\nu \in \{\mathsf{Rain}, \mathsf{Sunny}, \mathsf{Overcast}\}} \frac{|S_{\nu}|}{|S|} \mathsf{Entropy}(S_{\nu})$$

$$Gain(S, Outlook) = Entropy(S) - \frac{5}{14} Entropy(S_{Sunny}) - \frac{4}{14} Entropy(S_{Overca})$$

$$= 0.940 - \frac{5}{14} \times 0.971 - \frac{4}{14} \times 0 - \frac{5}{14} \times 0.971 = 0.940 - 0.347 - 0 - 0.347 = 0.2433 + 0.0000 + 0.0000$$

Information Gain



Learnt Decision Tree



Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

Calling ID3 on Outlook=Sunny

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

 Gain(S_{Outlook=Sunny}, Temp) = Entropy(2 Yes, 3 No) -(2/5)*Entropy(0 Yes, 2 No) -(2/5)*Entropy(1 Yes, 1 No) -(1/5)*Entropy(1 Yes, 0 No)
Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

 Gain(S_{Outlook=Sunny}, Temp) = Entropy(2 Yes, 3 No) -(2/5)*Entropy(0 Yes, 2 No) -(2/5)*Entropy(1 Yes, 1 No) -(1/5)*Entropy(1 Yes, 0 No)

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

- Gain(S_{Outlook=Sunny}, Temp) = Entropy(2 Yes, 3 No) -(2/5)*Entropy(0 Yes, 2 No) -(2/5)*Entropy(1 Yes, 1 No) -(1/5)*Entropy(1 Yes, 0 No)
- Gain(S_{Outlook=Sunny}, Humidity) = Entropy(2 Yes, 3 No) (2/5)*Entropy(2 Yes, 0 No) -(3/5)*Entropy(0 Yes, 3 No) ⇒ maximum possible for the set

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

- Gain(S_{Outlook=Sunny}, Temp) = Entropy(2 Yes, 3 No) -(2/5)*Entropy(0 Yes, 2 No) -(2/5)*Entropy(1 Yes, 1 No) -(1/5)*Entropy(1 Yes, 0 No)
- Gain(S_{Outlook=Sunny}, Humidity) = Entropy(2 Yes, 3 No) (2/5)*Entropy(2 Yes, 0 No) -(3/5)*Entropy(0 Yes, 3 No) ⇒ maximum possible for the set

Day	Temp	Humidity	Windy	Play
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

- Gain(S_{Outlook=Sunny}, Temp) = Entropy(2 Yes, 3 No) -(2/5)*Entropy(0 Yes, 2 No) -(2/5)*Entropy(1 Yes, 1 No) -(1/5)*Entropy(1 Yes, 0 No)
- Gain(S_{Outlook=Sunny}, Humidity) = Entropy(2 Yes, 3 No) (2/5)*Entropy(2 Yes, 0 No) -(3/5)*Entropy(0 Yes, 3 No) ⇒ maximum possible for the set
- Gain(S_{Outlook=Sunny}, Windy) = Entropy(2 Yes, 3 No) (3/5)*Entropy(1 Yes, 2 No) -(2/5)*Entropy(1 Yes, 1 No)

Learnt Decision Tree



Day	Temp	Humidity	Windy	Play
D4	Mild	High	Weak	Yes
D5	Cool	Normal	Weak	Yes
D6	Cool	Normal	Strong	No
D10	Mild	Normal	Weak	Yes
D14	Mild	High	Strong	No

• The attribute Windy gives the highest information gain

Learnt Decision Tree



Prediction for Decision Tree

Just walk down the tree!



Prediction for Decision Tree

Just walk down the tree!



Prediction for <High Humidity, Strong Wind, Sunny Outlook, Hot Temp> is ?

Prediction for Decision Tree

Just walk down the tree!



Prediction for <High Humidity, Strong Wind, Sunny Outlook, Hot Temp> is ? No

Apply the same rules, except when depth limit is reached, the leaf node is assigned the most common occurring value in that path.

Apply the same rules, except when depth limit is reached, the leaf node is assigned the most common occurring value in that path.

What is depth-0 tree (no decision) for the examples?

Apply the same rules, except when depth limit is reached, the leaf node is assigned the most common occurring value in that path.

What is depth-0 tree (no decision) for the examples? Always predicting Yes

Apply the same rules, except when depth limit is reached, the leaf node is assigned the most common occurring value in that path.

What is depth-0 tree (no decision) for the examples? Always predicting Yes

What is depth-1 tree (no decision) for the examples?

Apply the same rules, except when depth limit is reached, the leaf node is assigned the most common occurring value in that path.

What is depth-0 tree (no decision) for the examples? Always predicting Yes

What is depth-1 tree (no decision) for the examples?



Quick Question!

In the tennis dataset, why did "Outlook" have the highest information gain?

A) It was the first feature in the dataset

Quick Question!

- A) It was the first feature in the dataset
- B) When Outlook=Overcast, all examples have Play=Yes (pure subset)

Quick Question!

- A) It was the first feature in the dataset
- B) When Outlook=Overcast, all examples have Play=Yes (pure subset)
- C) It has the most possible values

Quick Question!

- A) It was the first feature in the dataset
- B) When Outlook=Overcast, all examples have Play=Yes (pure subset)
- C) It has the most possible values
- D) It was chosen randomly

Quick Question!

- A) It was the first feature in the dataset
- B) When Outlook=Overcast, all examples have Play=Yes (pure subset)
- C) It has the most possible values
- D) It was chosen randomly

Quick Question!

In the tennis dataset, why did "Outlook" have the highest information gain?

- A) It was the first feature in the dataset
- B) When Outlook=Overcast, all examples have Play=Yes (pure subset)
- C) It has the most possible values
- D) It was chosen randomly

Answer: B) When Outlook=Overcast, all examples have Play=Yes - This creates a pure subset with entropy=0, maximizing information gain.

Discrete Input, Real Output

Modified Dataset

Day	Outlook	Temp	Humidity	Wind	Minutes Played
D1	Sunny	Hot	High	Weak	20
D2	Sunny	Hot	High	Strong	24
D3	Overcast	Hot	High	Weak	40
D4	Rain	Mild	High	Weak	50
D5	Rain	Cool	Normal	Weak	60
D6	Rain	Cool	Normal	Strong	10
D7	Overcast	Cool	Normal	Strong	4
D8	Sunny	Mild	High	Weak	10
D9	Sunny	Cool	Normal	Weak	60
D10	Rain	Mild	Normal	Weak	40
D11	Sunny	Mild	High	Strong	45
D12	Overcast	Mild	High	Strong	40
D13	Overcast	Hot	Normal	Weak	35
D14	Rain	Mild	High	Strong	20

• Any guesses?

• Any guesses?

- Any guesses?
- Mean Squared Error

- Any guesses?
- Mean Squared Error

- Any guesses?
- Mean Squared Error
- MSE(S) = 311.34

- Any guesses?
- Mean Squared Error
- MSE(S) = 311.34

- Any guesses?
- Mean Squared Error
- MSE(S) = 311.34
- What about splitting criterion for regression?

- Any guesses?
- Mean Squared Error
- MSE(S) = 311.34
- What about splitting criterion for regression?

- Any guesses?
- Mean Squared Error
- MSE(S) = 311.34
- What about splitting criterion for regression?
- MSE Reduction (not Information Gain!)

- Any guesses?
- Mean Squared Error
- MSE(S) = 311.34
- What about splitting criterion for regression?
- MSE Reduction (not Information Gain!)
Measure of Impurity for Regression?

- Any guesses?
- Mean Squared Error
- MSE(S) = 311.34
- What about splitting criterion for regression?
- MSE Reduction (not Information Gain!)
- MSE Reduction = MSE(S) $\sum_{v} \frac{|S_v|}{|S|} MSE(S_v)$

Gain by splitting on Wind

Gain by splitting on Wind

Wind	Minutes Played
Weak	20
Strong	24
Weak	40
Weak	50
Weak	60
Strong	10
Strong	4
Weak	10
Weak	60
Weak	40
Strong	45
Strong	40
Weak	35
Strong	20

MSE(S)=311.34

Gain by splitting on Wind

Wind	Minutes Played
Weak	20
Strong	24
Weak	40
Weak	50
Weak	60
Strong	10
Strong	4
Weak	10
Weak	60
Weak	40
Strong	45
Strong	40
Weak	35
Strong	20

MSE(S)=311.34

Wind	Minutes Playe
Weak	20
Weak	40
Weak	50
Weak	60
Weak	10
Weak	60
Weak	40
Weak	35

$\mathsf{MSE}(S_{\mathsf{Wind}=\mathsf{Weak}}) = 277, \mathsf{Weight} = \frac{8}{14}$

Wind	Minutes Playe
Strong	24
Strong	10
Strong	4
Strong	45
Strong	40
Strong	20

 $MSE(S_{Wind=Strong}) = 218$, Weight = $\frac{6}{14}$

Correct calculation for Wind split:

 $\mathsf{MSE}\ \mathsf{Reduction} = \mathsf{MSE}(S) - \mathsf{Weighted}\ \mathsf{Average}\ \mathsf{MSE}$

$$= 311.34 - \left[\frac{8}{14} \times 277 + \frac{6}{14} \times 218\right] = 311.34 - [158.857 + 93.429] = 311.34$$

Key insight: MSE Reduction > 0 means the split improves our model!

For regression: Use MSE Reduction, NOT Information Gain!

For regression trees, what criterion do we use instead of Information Gain?

A) Information Gain

- A) Information Gain
- B) Gini Impurity

- A) Information Gain
- B) Gini Impurity
- C) Mean Squared Error (MSE) Reduction

- A) Information Gain
- B) Gini Impurity
- C) Mean Squared Error (MSE) Reduction
- D) Accuracy

- A) Information Gain
- B) Gini Impurity
- C) Mean Squared Error (MSE) Reduction
- D) Accuracy

For regression trees, what criterion do we use instead of Information Gain?

- A) Information Gain
- B) Gini Impurity
- C) Mean Squared Error (MSE) Reduction
- D) Accuracy

Answer: C) Mean Squared Error (MSE) Reduction - For regression, we minimize MSE instead of maximizing information gain.

MSE Reduction for Regression Trees



Learnt Tree



Learnt Tree



Real Input Discrete Output

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

• How do you find splits?

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).
- For the above example, we have 5 potential splits: 44, 54, 66, 76, 85

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).
- For the above example, we have 5 potential splits: 44, 54, 66, 76, 85
- Calculate the weighted impurity for each split

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- How do you find splits?
- Sort by attribute
- Find potential split points (midpoints).
- For the above example, we have 5 potential splits: 44, 54, 66, 76, 85
- Calculate the weighted impurity for each split
- Choose the split with the lowest impurity

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No
1		

• Consider split at 44

Temperature	PlayTennis
40	No
48	No
60	Yes
72	Yes
80	Yes
90	No
	Temperature 40 48 60 72 80 90

- Consider split at 44
- LHS has 1 No and 0 Yes; RHS has 3 Yes and 2 No

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 44
- LHS has 1 No and 0 Yes; RHS has 3 Yes and 2 No
- Entropy for LHS = 0, Entropy for RHS = 0.971

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 44
- LHS has 1 No and 0 Yes; RHS has 3 Yes and 2 No
- Entropy for LHS = 0, Entropy for RHS = 0.971
- Weighted Entropy = 0.971*5/6 = 0.808

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

• Consider split at 54

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 54
- LHS has 2 No and 0 Yes; RHS has 3 Yes and 1 No

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 54
- LHS has 2 No and 0 Yes; RHS has 3 Yes and 1 No
- Entropy for LHS = 0, Entropy for RHS = 0.811

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 54
- LHS has 2 No and 0 Yes; RHS has 3 Yes and 1 No
- Entropy for LHS = 0, Entropy for RHS = 0.811
- Weighted Entropy = 0.811*4/6 = 0.541

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No
1		

• Consider split at 66

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 66
- LHS has 2 No and 1 Yes; RHS has 2 Yes and 1 No

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 66
- LHS has 2 No and 1 Yes; RHS has 2 Yes and 1 No
- Entropy for LHS = 0.918, Entropy for RHS = 0.918

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 66
- LHS has 2 No and 1 Yes; RHS has 2 Yes and 1 No
- Entropy for LHS = 0.918, Entropy for RHS = 0.918
- Weighted Entropy = 0.918*3/6 + 0.918*3/6 = 0.918

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No
1		

• Consider split at 76

Temperature	PlayTennis
40	No
48	No
60	Yes
72	Yes
80	Yes
90	No
	Temperature 40 48 60 72 80 90

- Consider split at 76
- LHS has 2 No and 2 Yes; RHS has 1 Yes and 1 No

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 76
- LHS has 2 No and 2 Yes; RHS has 1 Yes and 1 No
- Entropy for LHS = 1, Entropy for RHS = 1
Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

- Consider split at 76
- LHS has 2 No and 2 Yes; RHS has 1 Yes and 1 No
- Entropy for LHS = 1, Entropy for RHS = 1
- Weighted Entropy = 1*4/6 + 1*2/6 = 1

Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

Notebook: decision-tree-real-input-discrete-output.html



Finding splits

Day	Temperature	PlayTennis
D1	40	No
D2	48	No
D3	60	Yes
D4	72	Yes
D5	80	Yes
D6	90	No

Notebook: decision-tree-real-input-discrete-output.html



Example (DT of depth 1)



Example (DT of depth 2)



Example (DT of depth 3)



Example (DT of depth 4)



Example (DT of depth 5)



Example (DT of depth 6)



Example (DT of depth 7)



Example (DT of depth 8)



Example (DT of depth 9)



Quick Question!

When finding splits for continuous features, how do we determine candidate split points?

A) Use all feature values as split points

Quick Question!

- A) Use all feature values as split points
- B) Use midpoints between consecutive sorted feature values

Quick Question!

- A) Use all feature values as split points
- B) Use midpoints between consecutive sorted feature values
- C) Use random values within the feature range

Quick Question!

- A) Use all feature values as split points
- B) Use midpoints between consecutive sorted feature values
- C) Use random values within the feature range
- D) Use only the minimum and maximum values

Quick Question!

- A) Use all feature values as split points
- B) Use midpoints between consecutive sorted feature values
- C) Use random values within the feature range
- D) Use only the minimum and maximum values

Quick Question!

When finding splits for continuous features, how do we determine candidate split points?

- A) Use all feature values as split points
- B) Use midpoints between consecutive sorted feature values
- C) Use random values within the feature range
- D) Use only the minimum and maximum values

Answer: B) Use midpoints between consecutive sorted feature values - This ensures we test all meaningful boundaries between different class regions. **Real Input Real Output**

Let us consider the dataset given below



57 / 87

What would be the prediction for decision tree with depth 0?



Prediction for decision tree with depth 0. Horizontal dashed line shows the predicted Y value. It is the average of Y values of all datapoints.



What would be the decision tree with depth 1?



Decision tree with depth 1



The Decision Boundary



What would be the decision tree with depth 2 ?



63 / 87

Decision tree with depth 2



The Decision Boundary



Feature is denoted by X and target by Y. Let the split be at X = s. Define regions: $R_1 = \{x : x \le s\}$ and $R_2 = \{x : x > s\}$. Feature is denoted by X and target by Y. Let the split be at X = s. Define regions: $R_1 = \{x : x \le s\}$ and $R_2 = \{x : x > s\}$.

For each region, compute the mean prediction:

$$c_{1} = \frac{1}{|R_{1}|} \sum_{x_{i} \in R_{1}} y_{i}$$

$$c_{2} = \frac{1}{|R_{2}|} \sum_{x_{i} \in R_{2}} y_{i}$$

Objective Function for Regression Trees

Feature is denoted by X and target by Y. Let the split be at X = s. Define regions: $R_1 = \{x : x \le s\}$ and $R_2 = \{x : x > s\}$.

For each region, compute the mean prediction:

$$c_1 = \frac{1}{|R_1|} \sum_{x_i \in R_1} y_i$$

$$c_2 = \frac{1}{|R_2|} \sum_{x_i \in R_2} y_i$$

The loss function is:

$$\mathsf{Loss}(s) = \sum_{x_i \in R_1} (y_i - c_1)^2 + \sum_{x_i \in R_2} (y_i - c_2)^2$$

Objective Function for Regression Trees

Feature is denoted by X and target by Y. Let the split be at X = s. Define regions: $R_1 = \{x : x \le s\}$ and $R_2 = \{x : x > s\}$.

For each region, compute the mean prediction:

$$c_1 = \frac{1}{|R_1|} \sum_{x_i \in R_1} y_i$$

$$c_2 = \frac{1}{|R_2|} \sum_{x_i \in R_2} y_i$$

The loss function is:

$$ext{Loss}(s) = \sum_{x_i \in R_1} (y_i - c_1)^2 + \sum_{x_i \in R_2} (y_i - c_2)^2$$

Our objective is to find the optimal split:

Algorithm: Finding the Optimal Split

1. Sort all data points (x_i, y_i) in increasing order of x_i .

- 1. Sort all data points (x_i, y_i) in increasing order of x_i .
- 2. Evaluate the loss function for all candidate splits:

$$s = \frac{x_i + x_{i+1}}{2}$$
 for $i = 1, 2, ..., n-1$

3. Select the split s^* that minimizes the loss function.

A Question!

Draw a regression tree for Y = sin(X), $0 \le X \le 2\pi$

A Question!

Dataset of Y = sin(X), $0 \le X \le 7$ with 10,000 points


A Question!

Regression tree of depth 1



A Question!

Decision Boundary



A Question!

Regression tree with no depth limit is too big to fit in a slide. It has of depth 4. The decision boundaries are in figure below.



What is the prediction function for a regression tree leaf node?

A) The median of target values in that region

- A) The median of target values in that region
- B) The mode of target values in that region

- A) The median of target values in that region
- B) The mode of target values in that region
- C) The mean of target values in that region

- A) The median of target values in that region
- B) The mode of target values in that region
- C) The mean of target values in that region
- D) A linear function of the features

- A) The median of target values in that region
- B) The mode of target values in that region
- C) The mean of target values in that region
- D) A linear function of the features

What is the prediction function for a regression tree leaf node?

- A) The median of target values in that region
- B) The mode of target values in that region
- C) The mean of target values in that region
- D) A linear function of the features

Answer: C) The mean of target values in that region - Each leaf predicts the average target value of training samples that reach that leaf.

Pruning and Overfitting

• Unpruned trees: Can grow very deep and complex

- Unpruned trees: Can grow very deep and complex
- **Perfect training accuracy**: Each leaf contains single training example

- Unpruned trees: Can grow very deep and complex
- **Perfect training accuracy**: Each leaf contains single training example
- But: Poor generalization to new data

- Unpruned trees: Can grow very deep and complex
- **Perfect training accuracy**: Each leaf contains single training example
- But: Poor generalization to new data
- Symptoms:

- Unpruned trees: Can grow very deep and complex
- **Perfect training accuracy**: Each leaf contains single training example
- But: Poor generalization to new data
- Symptoms:
 - High training accuracy, low test accuracy

- Unpruned trees: Can grow very deep and complex
- **Perfect training accuracy**: Each leaf contains single training example
- But: Poor generalization to new data
- Symptoms:
 - High training accuracy, low test accuracy
 - Very deep trees with many leaves

- Unpruned trees: Can grow very deep and complex
- **Perfect training accuracy**: Each leaf contains single training example
- But: Poor generalization to new data
- Symptoms:
 - High training accuracy, low test accuracy
 - Very deep trees with many leaves
 - Rules that are too specific to training data

- Unpruned trees: Can grow very deep and complex
- **Perfect training accuracy**: Each leaf contains single training example
- But: Poor generalization to new data
- Symptoms:
 - High training accuracy, low test accuracy
 - Very deep trees with many leaves
 - Rules that are too specific to training data
- Solution: Pruning to control model complexity

Stop growing tree before it becomes too complex:

• Maximum depth: Limit tree depth (e.g., max_depth = 5)

Stop growing tree before it becomes too complex:

- Maximum depth: Limit tree depth (e.g., max_depth = 5)
- Minimum samples per split: Don't split if node has i N samples

Stop growing tree before it becomes too complex:

- Maximum depth: Limit tree depth (e.g., max_depth = 5)
- Minimum samples per split: Don't split if node has i N samples
- Minimum samples per leaf: Ensure each leaf has
 M samples

Stop growing tree before it becomes too complex:

- Maximum depth: Limit tree depth (e.g., max_depth = 5)
- Minimum samples per split: Don't split if node has i N samples
- Minimum samples per leaf: Ensure each leaf has ≥ M samples
- Maximum features: Consider only subset of features at each split

Stop growing tree before it becomes too complex:

- Maximum depth: Limit tree depth (e.g., max_depth = 5)
- Minimum samples per split: Don't split if node has i N samples
- Minimum samples per leaf: Ensure each leaf has ≥ M samples
- Maximum features: Consider only subset of features at each split
- Minimum impurity decrease: Only split if improvement *i* threshold

Grow full tree, then remove unnecessary branches:

Grow full tree, then remove unnecessary branches:

• Algorithm:

1. Grow complete tree on training data

Grow full tree, then remove unnecessary branches:

- 1. Grow complete tree on training data
- 2. Use validation set to evaluate subtree performance

Grow full tree, then remove unnecessary branches:

- 1. Grow complete tree on training data
- 2. Use validation set to evaluate subtree performance
- 3. Remove branches that don't improve validation accuracy

Grow full tree, then remove unnecessary branches:

- 1. Grow complete tree on training data
- 2. Use validation set to evaluate subtree performance
- 3. Remove branches that don't improve validation accuracy
- 4. Repeat until no beneficial removals remain

Grow full tree, then remove unnecessary branches:

- 1. Grow complete tree on training data
- 2. Use validation set to evaluate subtree performance
- 3. Remove branches that don't improve validation accuracy
- 4. Repeat until no beneficial removals remain
- Cost Complexity Pruning: Minimize Error $+ \alpha \times$ Tree Size

Grow full tree, then remove unnecessary branches:

- 1. Grow complete tree on training data
- 2. Use validation set to evaluate subtree performance
- 3. Remove branches that don't improve validation accuracy
- 4. Repeat until no beneficial removals remain
- Cost Complexity Pruning: Minimize $\operatorname{Error} + \alpha \times \operatorname{Tree}$ Size
- Advantages: More thorough, can recover from early stopping mistakes

Grow full tree, then remove unnecessary branches:

- 1. Grow complete tree on training data
- 2. Use validation set to evaluate subtree performance
- 3. Remove branches that don't improve validation accuracy
- 4. Repeat until no beneficial removals remain
- Cost Complexity Pruning: Minimize $\operatorname{Error} + \alpha \times \operatorname{Tree} \operatorname{Size}$
- Advantages: More thorough, can recover from early stopping mistakes
- Disadvantages: More computationally expensive

Systematic approach to find optimal tree size:

• Cost function: $R_{\alpha}(T) = R(T) + \alpha |T|$

- Cost function: $R_{\alpha}(T) = R(T) + \alpha |T|$
 - R(T): Misclassification error on validation set

Cost Complexity Pruning Algorithm

- Cost function: $R_{\alpha}(T) = R(T) + \alpha |T|$
 - R(T): Misclassification error on validation set
 - |T|: Number of terminal nodes (tree size)

Cost Complexity Pruning Algorithm

- Cost function: $R_{\alpha}(T) = R(T) + \alpha |T|$
 - R(T): Misclassification error on validation set
 - |T|: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)

Cost Complexity Pruning Algorithm

- Cost function: $R_{\alpha}(T) = R(T) + \alpha |T|$
 - R(T): Misclassification error on validation set
 - |T|: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)
- Process:
- Cost function: $R_{\alpha}(T) = R(T) + \alpha |T|$
 - R(T): Misclassification error on validation set
 - |T|: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)
- Process:
 - 1. Start with full tree ($\alpha = 0$)

- Cost function: $R_{\alpha}(T) = R(T) + \alpha |T|$
 - R(T): Misclassification error on validation set
 - |T|: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)
- Process:
 - 1. Start with full tree ($\alpha = 0$)
 - 2. Gradually increase α

- Cost function: $R_{\alpha}(T) = R(T) + \alpha |T|$
 - R(T): Misclassification error on validation set
 - |T|: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)
- Process:
 - 1. Start with full tree ($\alpha = 0$)
 - 2. Gradually increase α
 - 3. At each α , prune branches that increase cost

- Cost function: $R_{\alpha}(T) = R(T) + \alpha |T|$
 - R(T): Misclassification error on validation set
 - |T|: Number of terminal nodes (tree size)
 - α : Complexity parameter (penalty for larger trees)
- Process:
 - 1. Start with full tree ($\alpha = 0$)
 - 2. Gradually increase α
 - 3. At each α , prune branches that increase cost
 - 4. Select α with best cross-validation performance

• Unpruned trees:

• Low bias (can fit complex patterns)

- Low bias (can fit complex patterns)
- High variance (sensitive to training data changes)

- Low bias (can fit complex patterns)
- High variance (sensitive to training data changes)
- Prone to overfitting

- Low bias (can fit complex patterns)
- High variance (sensitive to training data changes)
- Prone to overfitting
- Heavily pruned trees:

- Low bias (can fit complex patterns)
- High variance (sensitive to training data changes)
- Prone to overfitting
- Heavily pruned trees:
 - High bias (may miss important patterns)

- Low bias (can fit complex patterns)
- High variance (sensitive to training data changes)
- Prone to overfitting
- Heavily pruned trees:
 - High bias (may miss important patterns)
 - Low variance (more stable predictions)

- Low bias (can fit complex patterns)
- High variance (sensitive to training data changes)
- Prone to overfitting
- Heavily pruned trees:
 - High bias (may miss important patterns)
 - Low variance (more stable predictions)
 - Risk of underfitting

- Low bias (can fit complex patterns)
- High variance (sensitive to training data changes)
- Prone to overfitting
- Heavily pruned trees:
 - High bias (may miss important patterns)
 - Low variance (more stable predictions)
 - Risk of underfitting
- Optimal pruning: Balances bias and variance

- Low bias (can fit complex patterns)
- High variance (sensitive to training data changes)
- Prone to overfitting
- Heavily pruned trees:
 - High bias (may miss important patterns)
 - Low variance (more stable predictions)
 - Risk of underfitting
- Optimal pruning: Balances bias and variance
- Cross-validation: Essential for finding this balance

• Start simple: Begin with restrictive pre-pruning parameters

- Start simple: Begin with restrictive pre-pruning parameters
- Cross-validation: Always use CV to select pruning parameters

- Start simple: Begin with restrictive pre-pruning parameters
- Cross-validation: Always use CV to select pruning parameters
- Validation curves: Plot training/validation error vs. tree complexity

- Start simple: Begin with restrictive pre-pruning parameters
- Cross-validation: Always use CV to select pruning parameters
- Validation curves: Plot training/validation error vs. tree complexity
- Common parameters (sklearn):

- Start simple: Begin with restrictive pre-pruning parameters
- Cross-validation: Always use CV to select pruning parameters
- Validation curves: Plot training/validation error vs. tree complexity
- Common parameters (sklearn):
 - max_depth: Start with 3-10

- Start simple: Begin with restrictive pre-pruning parameters
- Cross-validation: Always use CV to select pruning parameters
- Validation curves: Plot training/validation error vs. tree complexity
- Common parameters (sklearn):
 - max_depth: Start with 3-10
 - min_samples_split: Try 10-100

- Start simple: Begin with restrictive pre-pruning parameters
- Cross-validation: Always use CV to select pruning parameters
- Validation curves: Plot training/validation error vs. tree complexity
- Common parameters (sklearn):
 - max_depth: Start with 3-10
 - min_samples_split: Try 10-100
 - min_samples_leaf: Try 5-50

- Start simple: Begin with restrictive pre-pruning parameters
- Cross-validation: Always use CV to select pruning parameters
- Validation curves: Plot training/validation error vs. tree complexity
- Common parameters (sklearn):
 - max_depth: Start with 3-10
 - min_samples_split: Try 10-100
 - min_samples_leaf: Try 5-50
 - ccp_alpha: Use for cost complexity pruning

- Start simple: Begin with restrictive pre-pruning parameters
- Cross-validation: Always use CV to select pruning parameters
- Validation curves: Plot training/validation error vs. tree complexity
- Common parameters (sklearn):
 - max_depth: Start with 3-10
 - min_samples_split: Try 10-100
 - min_samples_leaf: Try 5-50
 - ccp_alpha: Use for cost complexity pruning
- Domain knowledge: Consider interpretability requirements

Summary and Key Takeaways

• Interpretability an important goal

- Interpretability an important goal
- Decision trees: well known interpretable models

- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard

- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard
- Greedy approach:

- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard
- Greedy approach:
- Recursively split to maximize "performance gain"

- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard
- Greedy approach:
- Recursively split to maximize "performance gain"
- Issues:

- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard
- Greedy approach:
- Recursively split to maximize "performance gain"
- Issues:
 - Can overfit easily!

- Interpretability an important goal
- Decision trees: well known interpretable models
- Learning optimal tree is hard
- Greedy approach:
- Recursively split to maximize "performance gain"
- Issues:
 - Can overfit easily!
 - Empirically not as powerful as other methods







Weighted Entropy


Candidate Line: $X1 = 4(X1^*)$



Entropy of $X1 \leq X1^* = E_{\mathcal{S}(X1 < X1^*)}$

$$P(+) = \frac{0.1 + 0.1}{0.1 + 0.1 + 0.3} = \frac{2}{5}$$



Entropy of $X_1 > X_1^* = E_{S(X_1 > X_1^*)}$

 $P(+)=\frac{3}{5}$

