

# Decision Trees

Implementierungstechniken für  
Hauptspeicherdatenbanksysteme



Dominik Vinan

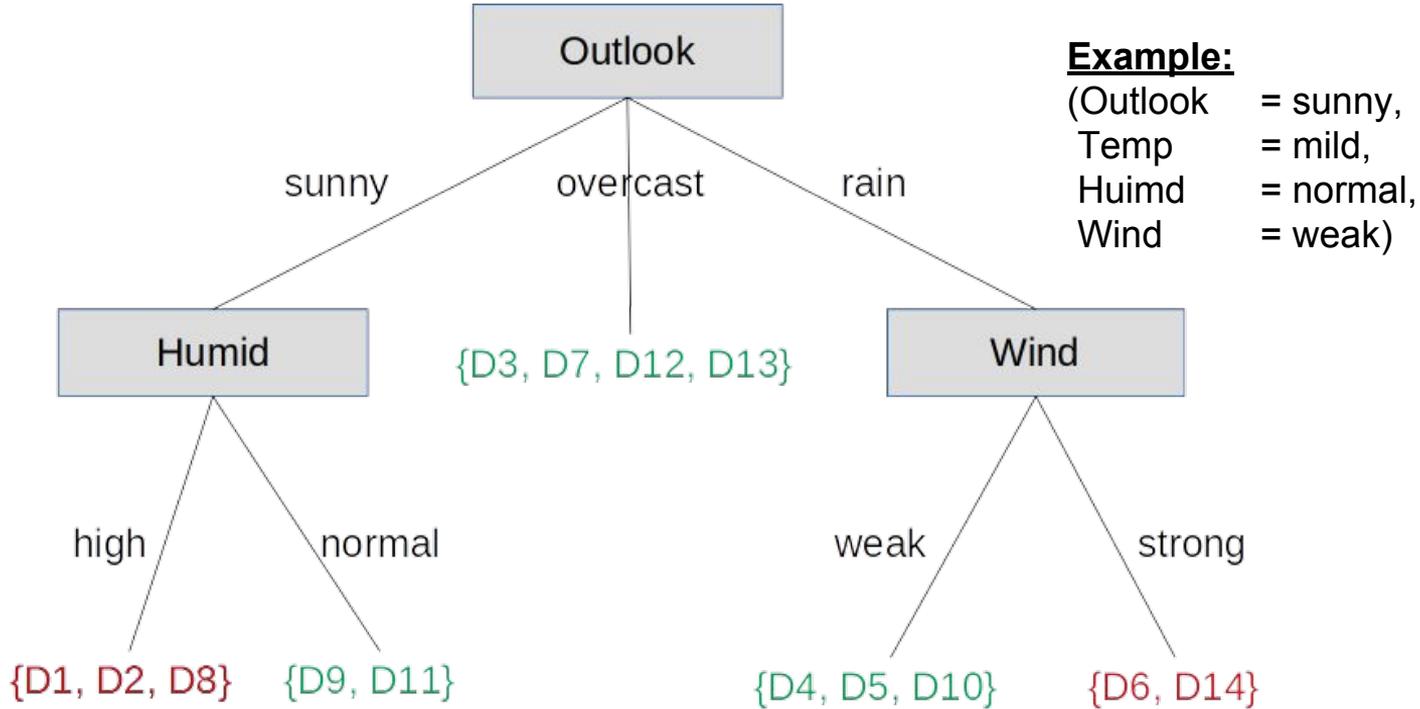
# Outline

- Introduction
- CART/TDIDT
  - ID3
  - C4.5
- Advantages & Disadvantages
- Implementation ID3

# Example Database

<u>Day</u>	<u>Outlook</u>	<u>Temp</u>	<u>Humid</u>	<u>Wind</u>	<u>Playtennis?</u>
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

# Introduction



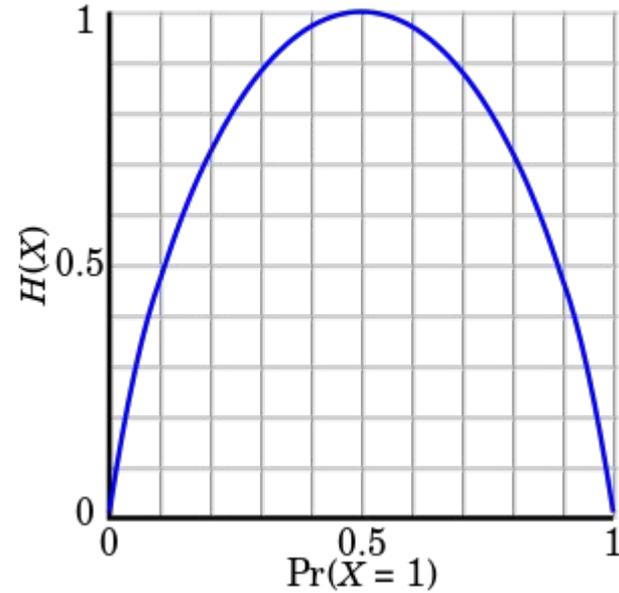
# CHAID & CART/TDIDT

- Top-Down Induction of Decision Trees (CART/TDIDT)
  - Non-Incremental approach i.e. needs to start over after change of training data.
  - Needs Pruning as Trees can become overly complex.
  - Examples: CART( Algorithm), ID3, C4.5, C5.0.
- Chi-square Automatic Interaction Detectors (CHAID)
  - Main Difference to CART:
    - Tree growth is limited => avoids pruning.

# Information

- Impurity (Entropy)
- Information Gain
  - Expected gain of information after splitting.

$$\text{gain}(A) = I(p, n) - H(A)$$



[https://en.wikipedia.org/wiki/Entropy\\_\(information\\_theory\)#/media/File:Binary\\_entropy\\_plot.svg](https://en.wikipedia.org/wiki/Entropy_(information_theory)#/media/File:Binary_entropy_plot.svg)

# Classification- vs. Regression Trees

## Classification Tree:

- Classifies categorical target values.
- E.g.: Response {'Yes', 'No'}.

## Regression Tree:

- Finds splitting value for continuous target values.
- E.g.: ExpectedTemp { 20.5, 10.7, 30.0, 17.7, ... }.

=> Predictors can be either numeric or categorical.

# Iterative Dichotomiser 3 (ID3)

- By J. Ross Quinlan
- Selection criterion: **Information Gain** or **Gain Ratio**.
- Information Gain
  - Bits of gained information.
  - Has a bias towards variables with multiple values.
- Gain Ratio
  - Takes number and size of branches into account.
  - Is not always defined.

# Example Database

<u>Day</u>	<u>Outlook</u>	<u>Temp</u>	<u>Humid</u>	<u>Wind</u>	<u>Playtennis?</u>
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

# ID3 Pseudocode

1. Select primary key, target attribute, and the dataset
2. If not (pure or stopping criterion met): Call ID3
  - a. Calculate Entropy & Information Gain for every attribute.
  - b. Select attribute X with  $\text{MAX}(\text{Information Gain})$ .
  - c. Make a tree node using attribute X.
  - d. Split dataset into subsets for every value of X.
  - e. Recursive call of ID3 for every subset.

# Example calculation ID3 (Root)

$$i(\text{Sunny}) = \frac{5}{14} \cdot \left( -\frac{2}{5} \cdot \log\left(\frac{2}{5}\right) - \frac{3}{5} \cdot \log\left(\frac{3}{5}\right) \right) = 0.346768$$

$$i(\text{Rain}) = \frac{5}{14} \cdot \left( -\frac{3}{5} \cdot \log\left(\frac{3}{5}\right) - \frac{2}{5} \cdot \log\left(\frac{2}{5}\right) \right) = 0.346768$$

$$i(\text{Overcast}) = \frac{4 \cdot (-1 \cdot \log(1))}{14} = 0$$

$$E(\text{Outlook}, D) = (i(\text{Sunny}) + i(\text{Overcast}) + i(\text{Rain})) = 0.693486$$

$$\text{gain}(\text{Outlook}) = I(D) - E(\text{Outlook}, D) = 0.940286 - 0.693486 = 0.2468$$

$$\text{gain}(\text{Humid}) = I(D) - E(\text{Humid}, 0) = 0.1518$$

$$\text{gain}(\text{Wind}) = I(D) - E(\text{Wind}, 0) = 0.0481$$

$$\text{gain}(\text{Temp}) = I(D) - E(\text{Temp}, 0) = 0.0292$$

# C4.5/C5.0

- Improvements over ID3
  - Handles continuous and categorical variables.
  - Handles missing values.
  - More efficient pruning.

=> C4.5 makes ID3 applicable in practice.

- C5.0 - Commercial implementation
  - Increased performance.
  - Less memory.
  - More precise.

# Advantages & Disadvantages

- + Human readable rules.
- + Limited computation power (for application).
- + Handles continuous and categorial values.
  
- Not optimal for predicting specific values.
- Growing and pruning trees is computationally complex.
- Inefficient for non-rectangular regions.

# Implementation

Day	Overcast	Temp	Humid	Wind	Day
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
...	...	...	...	...	...

**Thank You!**