# *** Applied Machine Learning Fundamentals ***
# Decision Trees and Ensembles

Daniel Wehner, M.Sc.

SAP SE / DHBW Mannheim

Winter term 2023/2024

# Lecture Overview

# Agenda for this Unit

❶ Introduction

❷ Iterative Dichotomizer (ID3)

❸ Extensions and Variants

❹ Ensemble Methods

❺ Wrap-Up

**Section:**

**Introduction**

What are Decision Trees?
An exemplary Tree
An alternative Tree

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
Ensemble Methods
Wrap-Up

**What are Decision Trees?**
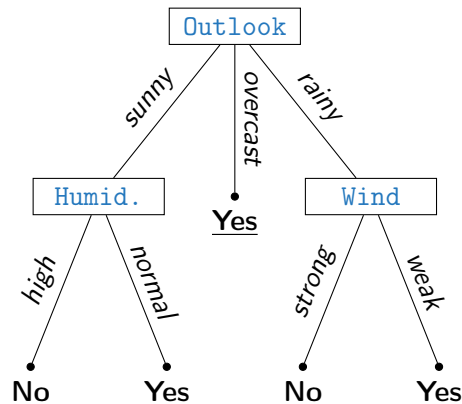An exemplary Tree
An alternative Tree

# What are Decision Trees?

- Decision trees are induced in a **supervised** fashion

- Originally invented by *Ross Quinlan* (1986)

- Decision trees are grown **recursively** → *'divide-and-conquer'*

- A decision tree consists of:

| | |
|---|---|
| **Nodes** | Each node corresponds to an attribute test |
| **Edges** | One edge per possible test outcome |
| **Leaves** | Class label to predict |

**Introduction**
Iterative Dichotomizer (ID3)
Extensions and Variants
Ensemble Methods
Wrap-Up

What are Decision Trees?
**An exemplary Tree**
An alternative Tree

## What we want...

| Outlook | Temperature | Humidity | Wind | PlayGolf |
|---------|-------------|----------|------|----------|
| sunny | hot | high | weak | **no** |
| sunny | hot | high | strong | **no** |
| overcast | hot | high | weak | **yes** |
| rainy | mild | high | weak | **yes** |
| rainy | cool | normal | weak | **yes** |
| rainy | cool | normal | strong | **no** |
| overcast | cool | normal | strong | **yes** |
| sunny | mild | high | weak | **no** |
| sunny | cool | normal | weak | **yes** |
| rainy | mild | normal | weak | **yes** |
| sunny | mild | normal | strong | **yes** |
| overcast | mild | high | strong | **yes** |
| overcast | hot | normal | weak | **yes** |
| rainy | mild | high | strong | **no** |
| rainy | mild | normal | strong | **???** |

**Introduction**
Iterative Dichotomizer (ID3)
Extensions and Variants
Ensemble Methods
Wrap-Up

What are Decision Trees?
**An exemplary Tree**
An alternative Tree

# Classifying new Instances

- Suppose we get a new instance:

  | Outlook | rainy |
  | Temperature | mild |
  | Humidity | normal |
  | Wind | strong |

- **What is its class?**

- Answer: **No**

**Introduction**
Iterative Dichotomizer (ID3)
Extensions and Variants
Ensemble Methods
Wrap-Up

What are Decision Trees?
An exemplary Tree
**An alternative Tree**

## Another Decision Tree...



Is this one better?

**Section:**
**Iterative Dichotomizer (ID3)**

Inductive Bias
Entropy and Information Gain
ID3 Algorithm

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

**Inductive Bias**
Entropy and Information Gain
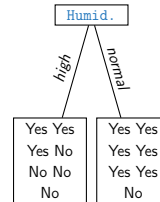ID3 Algorithm

# Inductive Bias of Decision Trees

- Complex models tend to **overfit** the data and **do not generalize well**

- Small decision trees are preferred

> **Occam's razor**:
> **'More things should not be used than are necessary.'**



- **Prefer the simplest hypothesis that fits the data!**

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
Entropy and Information Gain
ID3 Algorithm

# The Root of all Evil... Which Attribute to choose?

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
**Entropy and Information Gain**
ID3 Algorithm

# Finding a proper Attribute



- Simple and small trees are preferred
  - Data in successor node should be **as pure as possible**
  - I. e. nodes containing one class only are preferable
- **Question**: How can we express this thought as a mathematical formula?
- **Answer**:
  - Entropy (*Claude E. Shannon*)
  - Originates in the field of **information theory**

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
**Entropy and Information Gain**
ID3 Algorithm

# Measure of Impurity: Entropy

- Entropy is a measure of chaos in the data (measured in bits)

- **Example:** Consider two classes $A$ and $B$ ($\mathcal{C} = \{A, B\}$)

$$E(\{\boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}\}) \qquad \rightarrow 0 \qquad Bits$$
$$E(\{\boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}, B, B\}) \qquad \rightarrow 0.81 \qquad Bits$$
$$E(\{\boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}, \boldsymbol{A}, B, B, B, B\}) \qquad \rightarrow 1 \qquad Bit$$
$$E(\{\boldsymbol{A}, \boldsymbol{A}, B, B, B, B, B, B\}) \qquad \rightarrow 0.81 \qquad Bits$$
$$E(\{B, B, B, B, B, B, B, B\}) \qquad \rightarrow 0 \qquad Bits$$

> **If both classes are equally distributed, the entropy function $E$ reaches its maximum. Pure data sets have minimal entropy**.

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
**Entropy and Information Gain**
ID3 Algorithm

# Measure of Impurity: Entropy (Ctd.)



$$E(\mathcal{D}) = -\sum_{c \in \mathcal{C}} p_c \cdot \log_2 p_c$$

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
**Entropy and Information Gain**
ID3 Algorithm

# Measure of Impurity: Entropy (Ctd.)

**Entropy formula:**

$$E(\mathcal{D}) = -\sum_{c \in \mathcal{C}} p_c \cdot \log_2 p_c \tag{1}$$

- $p_c$ denotes the relative frequency of class $c \in \mathcal{C}$
- **Weather data:**

$$\mathcal{C} = \{yes, no\} \qquad \text{i. e.} \qquad p_{yes} = {}^9/_{14} \quad \text{and} \quad p_{no} = {}^5/_{14}$$

$$E(\mathcal{D}) = -\sum_{c \in \mathcal{C}} p_c \cdot \log_2 p_c = -({}^9/_{14} \cdot \log_2 {}^9/_{14} + {}^5/_{14} \cdot \log_2 {}^5/_{14}) = 0.9403$$

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
**Entropy and Information Gain**
ID3 Algorithm

# Quality of the Split: Average Entropy

- We still don't know which attribute to use for the split

- Calculate the entropy after each potential split

- **Average Entropy** after splitting by attribute `A`:

$$E(\mathcal{D}, \mathtt{A}) = \sum_{v \in \mathrm{dom}(\mathtt{A})} \frac{|\mathcal{D}_{\mathtt{A}=v}|}{|\mathcal{D}|} \cdot E(\mathcal{D}_{\mathtt{A}=v}) \qquad (2)$$

- Legend:

| | |
|---|---|
| `A` | Attribute |
| $\mathrm{dom}(\mathtt{A})$ | Possible values attribute `A` can take (domain of `A`) |
| $|\mathcal{D}_{\mathtt{A}=v}|$ | Number of examples satisfying $\mathtt{A} = v$ |

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
**Entropy and Information Gain**
ID3 Algorithm

# Quality of the Split: Average Entropy (Ctd.)

**Example:** Attribute `Outlook`

$$E(\mathcal{D}, \mathtt{Outlook}) = \sum_{v \in \mathsf{dom}(\mathtt{Outlook})} \frac{|\mathcal{D}_{\mathtt{Outlook}=v}|}{|\mathcal{D}|} \cdot E(\mathcal{D}_{\mathtt{Outlook}=v})$$

$$= \frac{5}{14} \cdot 0.9710 + \frac{5}{14} \cdot 0.9710 + \frac{4}{14} \cdot 0 \qquad = 0.6936$$

$$E(\mathcal{D}_{\mathtt{Outlook}=sunny}) = -(\frac{2}{5} \cdot \log_2(\frac{2}{5}) + \frac{3}{5} \cdot \log_2(\frac{3}{5})) \qquad = 0.9710$$

$$E(\mathcal{D}_{\mathtt{Outlook}=rainy}) = -(\frac{3}{5} \cdot \log_2(\frac{3}{5}) + \frac{2}{5} \cdot \log_2(\frac{2}{5})) \qquad = 0.9710$$

$$E(\mathcal{D}_{\mathtt{Outlook}=overcast}) = -(\frac{4}{4} \cdot \log_2(\frac{4}{4}) + \frac{0}{4} \cdot \log_2(\frac{0}{4})) \qquad = 0$$

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
**Entropy and Information Gain**
ID3 Algorithm

# Information Gain

- We have calculated the entropy before and after the split

- The difference of both is called the information gain (IG)

- Select the attribute with the highest IG

| Attribute | $E_{before}$ | $E_{after}$ | IG |
|---|---|---|---|
| `Outlook` | 0.9403 | 0.6936 | 0.2464 |
| `Temperature` | 0.9403 | 0.9111 | 0.0292 |
| `Humidity` | 0.9403 | 0.7885 | 0.1518 |
| `Wind` | 0.9403 | 0.8922 | 0.0481 |

- Attribute `Outlook` maximizes IG

- After the split: Remove attribute `Outlook`

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
**Entropy and Information Gain**
ID3 Algorithm

# Training Data after the Split by Attribute `Outlook`

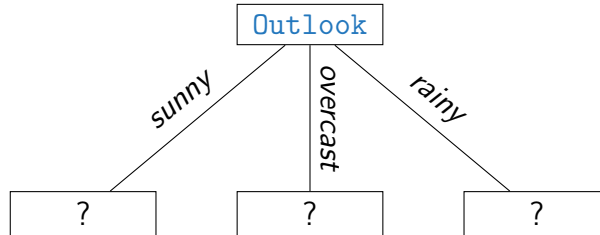| ~~Outlook~~ | Temperature | Humidity | Wind | PlayGolf |
|---|---|---|---|---|
| ~~sunny~~ | hot | high | weak | no |
| ~~sunny~~ | hot | high | strong | no |
| ~~sunny~~ | mild | high | weak | no |
| ~~sunny~~ | cool | normal | weak | yes |
| ~~sunny~~ | mild | normal | strong | yes |
| ~~rainy~~ | mild | high | weak | yes |
| ~~rainy~~ | cool | normal | weak | yes |
| ~~rainy~~ | cool | normal | strong | no |
| ~~rainy~~ | mild | normal | weak | yes |
| ~~rainy~~ | mild | high | strong | no |
| ~~overcast~~ | cool | normal | strong | yes |
| ~~overcast~~ | hot | high | weak | yes |
| ~~overcast~~ | mild | high | strong | yes |
| ~~overcast~~ | hot | normal | weak | yes |

- Data set $\mathcal{D}$ after the split
- We obtain three subsets (one per attribute value)
- Attribute `Outlook` is removed

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
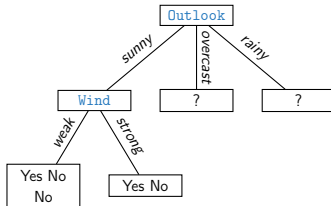Entropy and Information Gain
**ID3 Algorithm**

*Important*

# How to proceed?

- The algorithm is recursively applied to the resulting subsets
  1. Calculate entropy (before and after the split)
  2. Calculate information gain for each attribute
  3. Choose the attribute with max. information gain for the split
  4. In the current branch: Do not consider the attribute any more
  5. **Recursion** ↻ (`Go to 1`)

- Recursion stops as soon as the subset is pure

- In the example above the subset $\mathcal{D}_{\texttt{Outlook}=overcast}$ is already pure

- This algorithm is referred to as **ID3 (Iterative Dichotomizer)**

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
Entropy and Information Gain
**ID3 Algorithm**

# Step by Step: Construction of the Tree

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
Entropy and Information Gain
**ID3 Algorithm**

# Step by Step: Construction of the Tree (Ctd.)



- $IG(\texttt{Temperature}) = 0.571$

- $IG(\texttt{Humidity}) = \mathbf{0.971}$

- $IG(\texttt{Wind}) = 0.020$

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
Entropy and Information Gain
**ID3 Algorithm**

# Step by Step: Construction of the Tree (Ctd.)

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
Entropy and Information Gain
**ID3 Algorithm**

# Step by Step: Construction of the Tree (Ctd.)

Introduction
**Iterative Dichotomizer (ID3)**
Extensions and Variants
Ensemble Methods
Wrap-Up

Inductive Bias
Entropy and Information Gain
**ID3 Algorithm**

# Step by Step: Construction of the Tree (Ctd.)

---

**Algorithm 1:** ID3 Algorithm (Iterative Dichotomizer)

**Input:** Training set $\mathcal{D}$, Attribute list *Attr_List*

1 Create a node $N$
2 **if** *all tuples in $\mathcal{D}$ have class c* **then**
3      **return** $N$ *as leaf node labeled with class c* // subset is pure
4 **if** $|Attr\_List| = 0$ **then**
5      **return** $N$ *as leaf node labeled with majority class in $\mathcal{D}$* // no attribute left
6 Find best split attribute $\mathtt{A}^*$ and label node $N$ with $\mathtt{A}^*$
7 *Attr_List* $\longleftarrow$ *Attr_List* $\setminus \{\mathtt{A}^*\}$ // remove attribute from list
8 **forall** $v \in dom(\mathtt{A}^*)$ **do**
9      Let $\mathcal{D}_{\mathtt{A}^*=v}$ be the set of tuples in $\mathcal{D}$ that satisfy $\mathtt{A}^* = v$
10      **if** $|\mathcal{D}_{\mathtt{A}^*=v}| = 0$ **then**
11          Attach leaf labeled with majority class in $\mathcal{D}$ to node $N$
12      **else**
13          Attach node returned by *ID3*($\mathcal{D}_{\mathtt{A}^*=v}$, *Attr_List*) // recursion
14 **return** $N$

---

**Section:**

**Extensions and Variants**

Other Measures of Impurity
Highly-branching Attributes
Numeric Attributes
Regression Trees

Introduction
Iterative Dichotomizer (ID3)
**Extensions and Variants**
Ensemble Methods
Wrap-Up

Other Measures of Impurity
Highly-branching Attributes
Numeric Attributes
Regression Trees

# An Alternative to Information Gain: Gini Index

**Gini index:**

$$Gini(\mathcal{D}) = \sum_{c \in \mathcal{C}} p_c \cdot (1 - p_c) = 1 - \sum_{c \in \mathcal{C}} p_c^2 \qquad (3)$$

- Used e. g. in **CART (Classification and Regression Trees)**
- **Gini gain** could be defined analogously to IG
  *(usually not done)*

Introduction
Iterative Dichotomizer (ID3)
**Extensions and Variants**
Ensemble Methods
Wrap-Up

**Other Measures of Impurity**
Highly-branching Attributes
Numeric Attributes
Regression Trees

# Why not use the Error as a splitting Criterion?

- The bias towards pure leaves is **not strong enough**

- **Example:**

| | |
|---|---|
| 40 of A | 60 of A |
| 60 of A | 40 of B |

Split 1 (vertical axis), Split 2 (horizontal axis)

Error without splitting:
**20 %**

Error after splitting:
**20 %**

**Both splits don't improve the error.**
**But together they give a perfect split!**

Introduction
Iterative Dichotomizer (ID3)
**Extensions and Variants**
Ensemble Methods
Wrap-Up

**Other Measures of Impurity**
Highly-branching Attributes
Numeric Attributes
Regression Trees

# Summary: Impurity Measures

Introduction
Iterative Dichotomizer (ID3)
**Extensions and Variants**
Ensemble Methods
Wrap-Up

Other Measures of Impurity
**Highly-branching Attributes**
Numeric Attributes
Regression Trees

# Highly-branching Attributes

**Attributes with a large number of values are problematic, since the leaves are not 'backed' with sufficient data examples.**

**In extreme cases only one example per node (e. g. IDs)**

This may lead to:

- **Overfitting** (Selection of attributes which are not optimal for prediction)
- **Fragmentation** (Data is fragmented into (too) many small sets)

Introduction
Iterative Dichotomizer (ID3)
**Extensions and Variants**
Ensemble Methods
Wrap-Up

Other Measures of Impurity
**Highly-branching Attributes**
Numeric Attributes
Regression Trees

# Highly-branching Attributes (Ctd.)



- Entropy before was 0.9403, Entropy after split is 0

- IG($\mathcal{D}$, Day) = 0.9403

- Attribute Day would be chosen for the split ⇒ **Bad for prediction** ☠

Introduction
Iterative Dichotomizer (ID3)
**Extensions and Variants**
Ensemble Methods
Wrap-Up

Other Measures of Impurity
**Highly-branching Attributes**
Numeric Attributes
Regression Trees

# Highly-branching Attributes (Ctd.)

- Calculate the **intrinsic information (IntI)**:

$$IntI(\mathcal{D}, \mathtt{A}) = - \sum_{v \in \mathsf{dom}(\mathtt{A})} \frac{|\mathcal{D}_{\mathtt{A}=v}|}{|\mathcal{D}|} \cdot \log_2 \frac{|\mathcal{D}_{\mathtt{A}=v}|}{|\mathcal{D}|} \qquad (4)$$

- Attributes with high *IntI* are **less useful** (high fragmentation)
- New splitting heuristic **Gain ratio (GR)**:

$$GR(\mathcal{D}, \mathtt{A}) = \frac{IG(\mathcal{D}, \mathtt{A})}{IntI(\mathcal{D}, \mathtt{A})} \qquad (5)$$

Introduction
Iterative Dichotomizer (ID3)
**Extensions and Variants**
Ensemble Methods
Wrap-Up

Other Measures of Impurity
**Highly-branching Attributes**
Numeric Attributes
Regression Trees

# Highly-branching Attributes (Ctd.)

- Intrinsic information for attribute Day:

$$IntI(\mathcal{D}, \text{Day}) = 14 \cdot (-\frac{1}{14} \cdot \log_2(\frac{1}{14})) = 3.807 \tag{6}$$

- Gain ratio for attribute Day:

$$GR(\mathcal{D}, \text{Day}) = \frac{0.9403}{3.807} = 0.246 \tag{7}$$

In this case the attribute Day would still be chosen. Be careful what features to include into the training data set! **(feature engineering is important!)**

Introduction
Iterative Dichotomizer (ID3)
**Extensions and Variants**
Ensemble Methods
Wrap-Up

Other Measures of Impurity
Highly-branching Attributes
**Numeric Attributes**
Regression Trees

# Handling numeric Attributes

- Usually, only **binary splits** are considered, e. g.:
  - Temperature $< 48$
  - CPU $> 24$
  - **Not:** $24 \leqslant$ Temperature $\leqslant 31$

- To support multiple splits, the attribute is **not removed**

  *(the same attribute can be used again for another split)*

- **Problem:** There is an **infinite number** of possible splits!

- **Solution:** Discretize range (fixed step size, ...)

- **Splitting on numeric attributes is computationally demanding!**

Introduction
Iterative Dichotomizer (ID3)
**Extensions and Variants**
Ensemble Methods
Wrap-Up

Other Measures of Impurity
Highly-branching Attributes
**Numeric Attributes**
Regression Trees

## Handling numeric Attributes (Ctd.)

- Consider the attribute `Temperature`:
  Use **numerical values** instead of discrete values like *cool*, *mild*, *hot*:

| 64 | 65 | 68 | 69 | 70 | 71 | 72 | 72 | 75 | 75 | 80 | 81 | 83 | 85 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| Yes | No | Yes | Yes | Yes | No | No | Yes | Yes | Yes | No | Yes | Yes | No |

- `Temperature` $< 71.5$

  yes: 4 | no: 2

- `Temperature` $\geqslant 71.5$

  yes: 5 | no: 3

$$E(\mathcal{D}, \texttt{Temp.}) = {}^{6}\!/\!{}_{14} \cdot E(\texttt{Temp.} < 71.5) + {}^{8}\!/\!{}_{14} \cdot E(\texttt{Temp.} \geqslant 71.5) = 0.939$$

Introduction
Iterative Dichotomizer (ID3)
**Extensions and Variants**
Ensemble Methods
Wrap-Up

Other Measures of Impurity
Highly-branching Attributes
**Numeric Attributes**
Regression Trees

# Handling numeric Attributes (Ctd.)

| Sorted Values | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| No | No | No | Yes | Yes | Yes | No | No | No | No |
| **Taxable Income** | | | | | | | | | |
| 60 | 70 | 75 | 85 | 90 | 95 | 100 | 120 | 125 | 220 |

| Split | 55 | | 65 | | 72 | | 80 | | 87 | | 92 | | 97 | | 110 | | 122 | | 172 | | 230 | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ⩽ | > | ⩽ | > | ⩽ | > | ⩽ | > | ⩽ | > | ⩽ | > | ⩽ | > | ⩽ | > | ⩽ | > | ⩽ | > | ⩽ | > |
| Yes | 0 | 3 | 0 | 3 | 0 | 3 | 0 | 3 | 1 | 2 | 2 | 1 | 3 | 0 | 3 | 0 | 3 | 0 | 3 | 0 | 3 | 0 |
| No | 0 | 7 | 1 | 6 | 2 | 5 | 3 | 4 | 3 | 4 | 3 | 4 | 3 | 4 | 4 | 3 | 5 | 2 | 6 | 1 | 7 | 0 |
| Gini | 0.420 | | 0.400 | | 0.375 | | 0.343 | | 0.417 | | 0.400 | | <u>0.300</u> | | 0.343 | | 0.375 | | 0.400 | | 0.420 | |

Introduction
Iterative Dichotomizer (ID3)
**Extensions and Variants**
Ensemble Methods
Wrap-Up

Other Measures of Impurity
Highly-branching Attributes
Numeric Attributes
**Regression Trees**

# Regression Trees

- Prediction of continuous variables

- Predict average value of all examples in the leaf

- Split the data such that variance in the leaves is minimized

- **Termination criterion is important, otherwise single point per leaf!**

**Standard deviation reduction (SDR):**

$$SDR(\mathcal{D}, \mathtt{A}) = SD(\mathcal{D}) - \sum_{v \in dom(\mathtt{A})} \frac{|\mathcal{D}_{\mathtt{A}=v}|}{|\mathcal{D}|} \cdot SD(\mathcal{D}_{\mathtt{A}=v}) \qquad (8)$$

**Section:**

# Ensemble Methods

Introduction to Ensembles
Bagging and Randomization

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
**Ensemble Methods**
Wrap-Up

Introduction to Ensembles
Bagging and Randomization

# Introduction Ensemble Methods

- **Key Idea**: Don't learn a single classifier, but a **set of classifiers**
- Combine the predictions of the single classifiers to obtain the final prediction

**Problem:** How can we induce multiple classifiers from a single data set without getting the same classifier over and over again? **We want to have diverse classifiers, otherwise the ensemble is useless!**

- Basic techniques:
  - **Bagging**
  - **Boosting** (not covered)
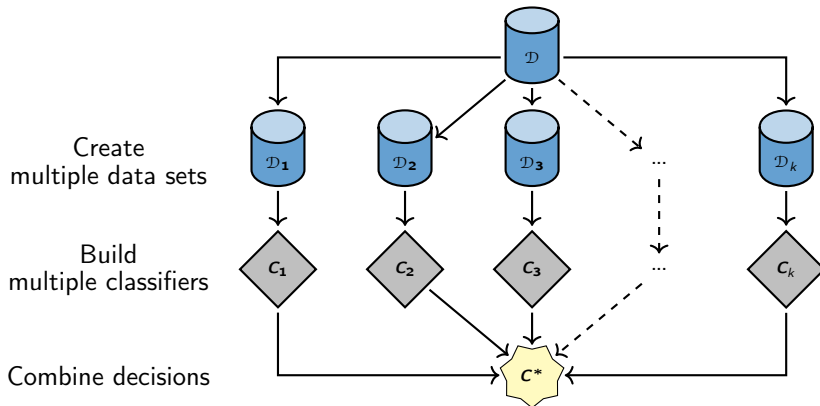  - **Stacking** (not covered)

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
**Ensemble Methods**
Wrap-Up

Introduction to Ensembles
Bagging and Randomization

# What is the Advantage?

- Consider the following:
  - There are 25 **independent** base classifiers
  - Independence assumption: Probability of misclassification **does not** depend on other classifiers in the ensemble
  - Usually, this assumption does not fully hold in practice
  - Each classifier has an error rate of $\varepsilon = 0.35$
- The ensemble makes a wrong prediction **if the majority is wrong** ($\Rightarrow$ i. e. at least 13)

$$\varepsilon_{ensemble} = \sum_{u=13}^{25} \binom{25}{u} \cdot \varepsilon^u \cdot (1-\varepsilon)^{25-u} \approx 0.06 \ll \varepsilon \qquad (9)$$

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
**Ensemble Methods**
Wrap-Up

Introduction to Ensembles
Bagging and Randomization

# Bagging: General Approach

Bagging $\widehat{=}$ **B**ootstrap **Agg**regat**ing**

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
**Ensemble Methods**
Wrap-Up

Introduction to Ensembles
Bagging and Randomization

# Creating the Bootstrap Samples

- How to generate multiple data sets which are different?
- **Solution**: Use sampling with replacement

| Original Data | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| **Bagging (Round 1)** | 7 | 8 | 10 | 8 | 2 | 5 | 10 | 10 | 5 | 9 |
| **Bagging (Round 2)** | 1 | 4 | 9 | 1 | 2 | 3 | 2 | 7 | 3 | 2 |
| **Bagging (Round 3)** | 1 | 8 | 5 | 10 | 5 | 5 | 9 | 6 | 3 | 7 |

- Some examples may appear **in more than one set**
- Some examples may appear **more than once** in one set
- Some examples may **not appear at all**

**Algorithm 2:** Bagging Algorithm

**Input:** Training set $\mathcal{D}$, number of base classifiers $k$

**1 Training:**

**2 forall** $u \in \{1, 2, \ldots, k\}$ **do**

3      Draw a bootstrap sample $\mathcal{D}_u$ with replacement from $\mathcal{D}$

4      Learn a classifier $C_u$ from $\mathcal{D}_u$

5      Add classifier $C_u$ to the ensemble

**6 Prediction:**

**7 forall** *unlabeled instances* **do**

8      Get predictions from all classifiers $C_u$

**9 return** *Class which receives the majority of votes (combined classifier $C^*$)*

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
Ensemble Methods
Wrap-Up

Introduction to Ensembles
Bagging and Randomization

# Bagging Variations

- The bootstrap samples had equal size and were drawn with replacement
- Also conceivable:
    1. **Varying the size** of the bootstrap samples
    2. Sampling **without replacement** $\Rightarrow$ **Pasting**
    3. **Sampling of features**, not instances
        - Not all features are available in all bootstrap samples
        - This is how **random forests** work
    4. Creating **heterogeneous ensembles**
       (neural networks, decision trees, support vector machines, ...)

Introduction
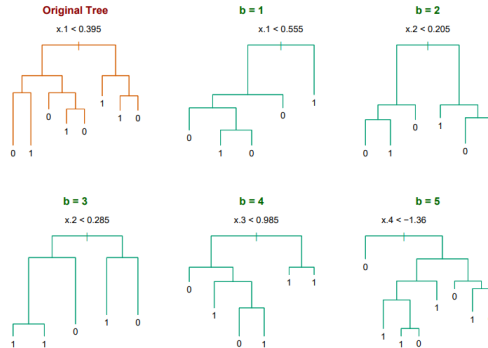Iterative Dichotomizer (ID3)
Extensions and Variants
**Ensemble Methods**
Wrap-Up

Introduction to Ensembles
Bagging and Randomization

# Bagged Decision Trees



Figure: Bagged decision trees; cf. Hastie 2008, page 284

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
**Ensemble Methods**
Wrap-Up

Introduction to Ensembles
Bagging and Randomization

# Randomization

- Why not **randomizing the algorithm** instead of the data?
- Some algorithms already do that: E. g. neural networks (random initialization of weights)
- Especially greedy algorithms can be randomized:
  - Pick from the options **randomly**, instead of picking the best one
  - E. g. decision trees: Do not choose attribute with highest information gain

**A random forest combines randomization and bagging.**

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
**Ensemble Methods**
Wrap-Up

Introduction to Ensembles
Bagging and Randomization

# Random Forest Algorithm

- Ensemble of decision trees
- Combines **bagging** and **random attribute subset selection**
- Build decision tree from a bootstrap sample
- Select best split attribute among a random subset of $f$ attributes

A random forest selects the best splitting attributes from the set of features available, but the globally best features **may not** be available.

## Algorithm 3: Random Forest Algorithm

**Input:** Training set $\mathcal{D}$, number of base classifiers $k$

1 **Training:**

2 **for** $u \in \{1, 2, \ldots, k\}$ **do**

3      Create a bootstrap sample from $\mathcal{D}$ (e. g. with replacement) $\Rightarrow$ **Bagging**

4      **begin**

5          Grow the tree

6          At every node: Randomly choose $f$ attributes to be considered for the split

7          $\Rightarrow$ **Randomization**

8          Do not prune tree $C_u$

9      Add tree $C_u$ to the ensemble

10 **Prediction:**

11 **forall** *unlabeled instances* **do**

12      Get predictions from all classifiers $C_u$

13 **return** *Class which receives the majority of votes (combined classifier $C^*$)*

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
**Ensemble Methods**
Wrap-Up

Introduction to Ensembles
Bagging and Randomization

# ExtraTrees (Randomization 2.0)

- One more step of randomization $\Rightarrow$ **Ext**remely **Ra**ndomized **Trees**

- The general approach is the same as for random forests
  **But:**
  - Instead of choosing the optimal split point...
  - ...it is selected randomly
  - The decision tree is grown without having to calculate entropy

- It is **much faster** (due to less computation)

The large number of classifiers compensates for suboptimal splits.

**Section:**

**Wrap-Up**

Summary
Self-Test Questions
Lecture Outlook

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
Ensemble Methods
**Wrap-Up**

**Summary**
Self-Test Questions
Lecture Outlook

# Summary

- **Decision trees:**
  - The construction of decision trees is guided by an **impurity measure**, e. g. entropy or Gini
  - Recursively select features which **maximize the information gain**
  - Decision trees can handle **numeric attributes** and **continuous output**
- **Ensembles:**
  - Usually, ensembles allow for a significant error reduction
  - **Bagging:** Sample diverse data sets from underlying data
  - **Random forests** combine bagging and randomization

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
Ensemble Methods
**Wrap-Up**

Summary
**Self-Test Questions**
Lecture Outlook

## Self-Test Questions

1. What is an inductive bias? What is the inductive bias of decision trees?

2. Explain what Occam's razor is.

3. What does entropy measure? How do you compute the information gain?

4. True or false? *'Pure data sets have maximal entropy.'*

5. What is the advantage of ensemble methods?

6. What is bagging?

7. Explain what a random forest does.

Introduction
Iterative Dichotomizer (ID3)
Extensions and Variants
Ensemble Methods
**Wrap-Up**

Summary
Self-Test Questions
**Lecture Outlook**

# What's next...?

| | |
|---|---|
| **Unit I** | Machine Learning Introduction |
| **Unit II** | Mathematical Foundations |
| **Unit III** | Bayesian Decision Theory |
| **Unit IV** | Regression |
| **Unit V** | Classification I |
| **Unit VI** | **Evaluation** |
| **Unit VII** | Classification II |
| **Unit VIII** | Clustering |
| **Unit IX** | Dimensionality Reduction |

# Thank you very much for the attention!

**Topic:**   *** Applied Machine Learning Fundamentals *** Decision Trees and Ensembles
**Term:**   Winter term 2023/2024

**Contact:**
Daniel Wehner, M.Sc.
SAP SE / DHBW Mannheim
daniel.wehner@sap.com

## Do you have any questions?