

*** Applied Machine Learning Fundamentals ***

Decision Trees and Ensembles

M. Sc. Daniel Wehner

SAP SE

Winter term 2019/2020



Find all slides on [GitHub](#)

Lecture Overview

Unit I	Machine Learning Introduction
Unit II	Mathematical Foundations
Unit III	Bayesian Decision Theory
Unit IV	Probability Density Estimation
Unit V	Regression
Unit VI	Classification I
Unit VII	Evaluation
Unit VIII	Classification II
Unit IX	Clustering
Unit X	Dimensionality Reduction

Agenda for this Unit

① Introduction

② Iterative Dichotomizer (ID3)

Inductive Bias
Entropy and Information Gain
ID3 Algorithm

③ Extensions and Variants

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

④ Ensemble Methods

Introduction to Ensembles
Bagging and Randomization
Boosting
Stacking
Error-correcting Output Codes

⑤ Wrap-Up

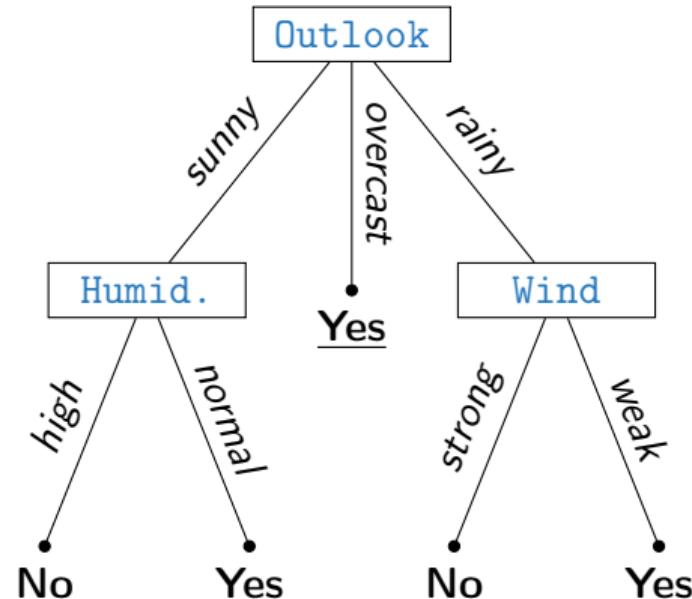
Summary
Self-Test Questions
Lecture Outlook
Recommended Literature and further Reading
Meme of the Day

Section: Introduction



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	???



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:

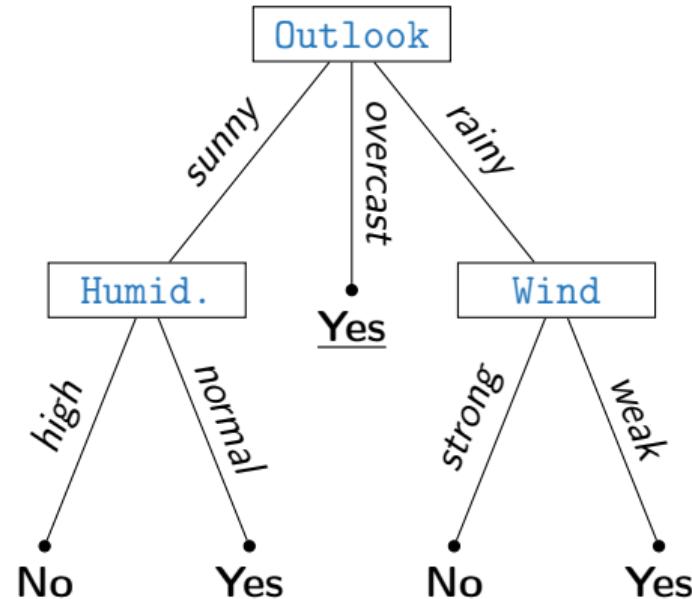
Nodes	Each node corresponds to an attribute test
Edges	One edge per possible test outcome
Leaves	Class label to predict

Classifying new Instances

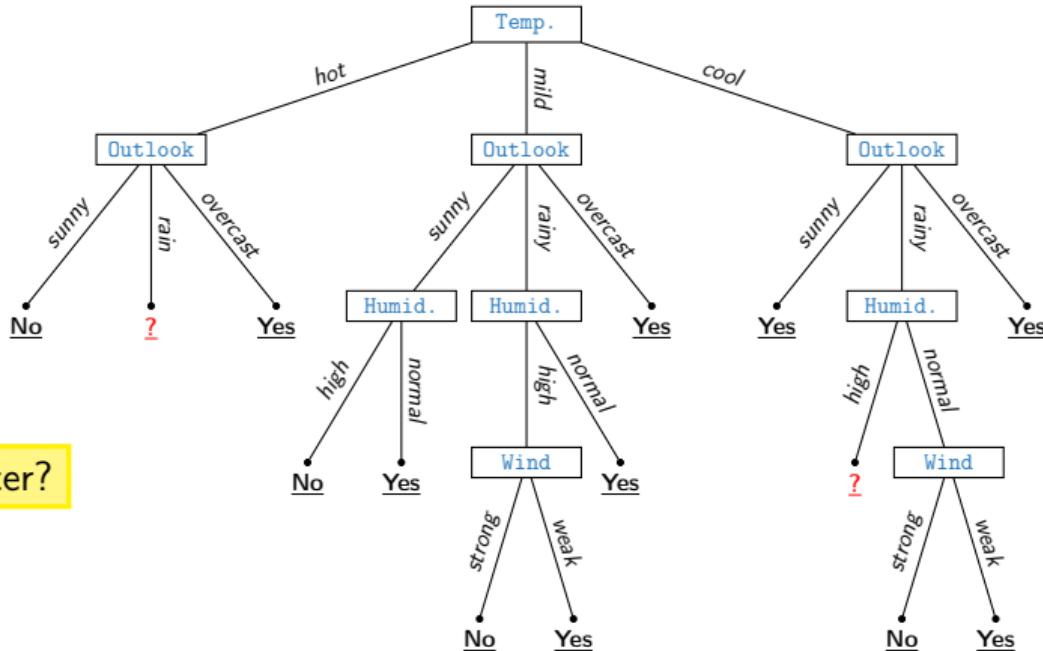
- Suppose we get a new instance:

Outlook	rainy
Temperature	mild
Humidity	normal
Wind	strong

- What is its class?
- Answer: No



Another Decision Tree...



Section:
Iterative Dichotomizer (ID3)



Inductive Bias of Decision Trees

- Complex models tend to **overfit** the data and **do not generalize well**
- Small decision trees are preferred

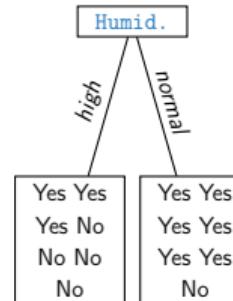
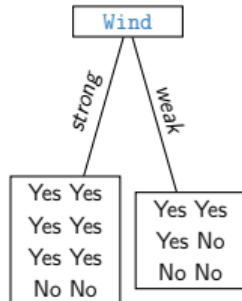
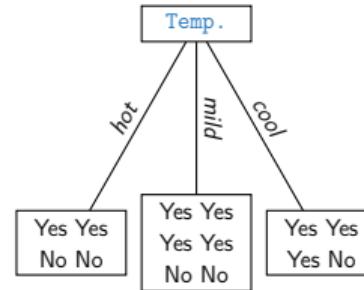
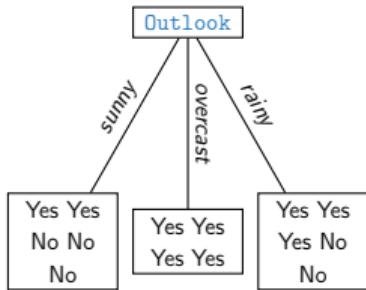
Occam's razor:

'More things should not be used than are necessary.'



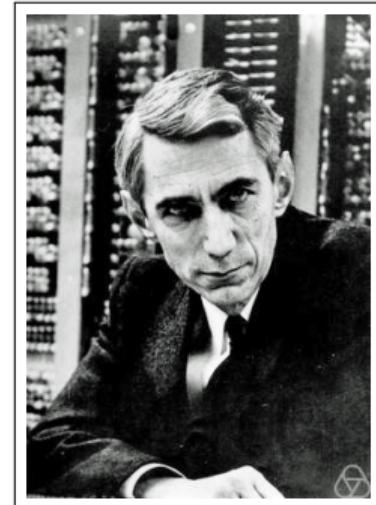
- **Prefer the simplest hypothesis that fits the data!**

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



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**



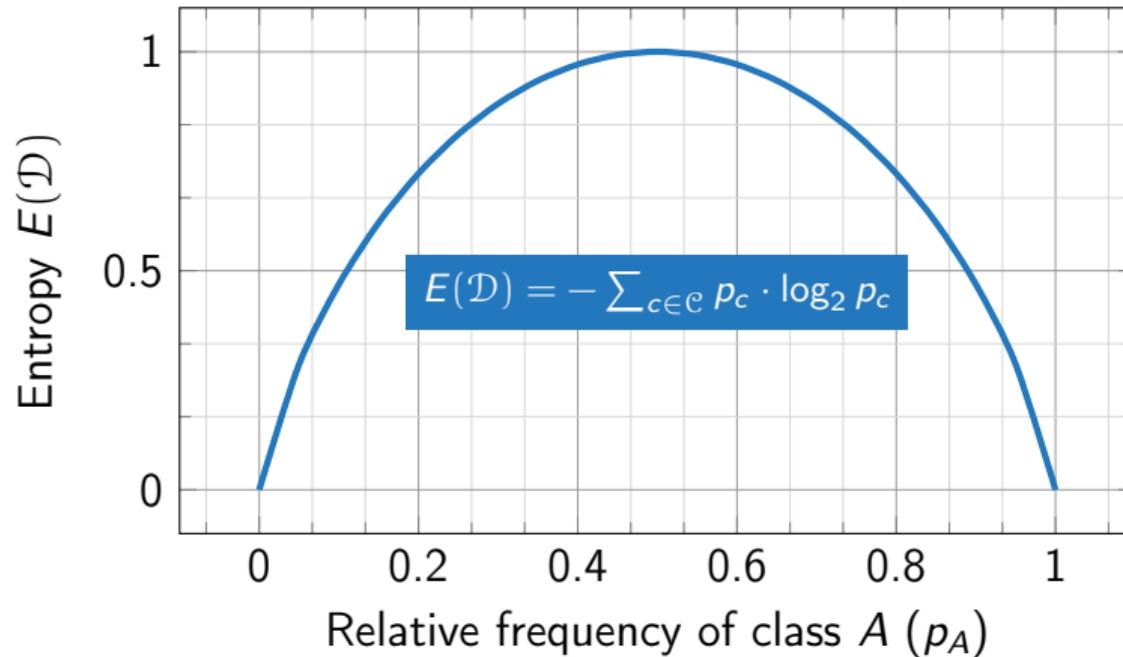
Measure of Impurity: Entropy

- Entropy is a measure of chaos in the data (measured in bits)
- **Example:** Consider two classes A and B ($\mathcal{C} = \{A, B\}$)

$$\begin{array}{lll} E(\{A, A, A, A, A, A, A, A\}) & \rightarrow 0 & \text{Bits} \\ E(\{A, A, A, A, A, A, B, B\}) & \rightarrow 0.81 & \text{Bits} \\ E(\{A, A, A, A, B, B, B, B\}) & \rightarrow 1 & \text{Bit} \\ E(\{A, A, B, B, B, B, B, B\}) & \rightarrow 0.81 & \text{Bits} \\ E(\{B, B, B, B, B, B, B, B\}) & \rightarrow 0 & \text{Bits} \end{array}$$

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

Measure of Impurity: Entropy (Ctd.)



Measure of Impurity: Entropy (Ctd.)

Entropy formula:

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

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

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

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

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:

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

- Legend:

A

Attribute

$\text{dom}(A)$

Possible values attribute A can take (domain of A)

$|\mathcal{D}_{A=v}|$

Number of examples satisfying $A = v$

Quality of the Split: Average Entropy (Ctd.)

Example: Attribute Outlook

$$E(\mathcal{D}, \text{Outlook}) = \sum_{v \in \text{dom}(\text{Outlook})} \frac{|\mathcal{D}_{\text{Outlook}=v}|}{|\mathcal{D}|} \cdot E(\mathcal{D}_{\text{Outlook}=v})$$
$$= \frac{5}{14} \cdot 0.9710 + \frac{5}{14} \cdot 0.9710 + \frac{4}{14} \cdot 0 = \mathbf{0.6936}$$

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

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

$$E(\mathcal{D}_{\text{Outlook}=\text{overcast}}) = -(4/4 \cdot \log_2(4/4) + 0/4 \cdot \log_2(0/4)) = 0$$

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

Training Data after the Split by Attribute Outlook

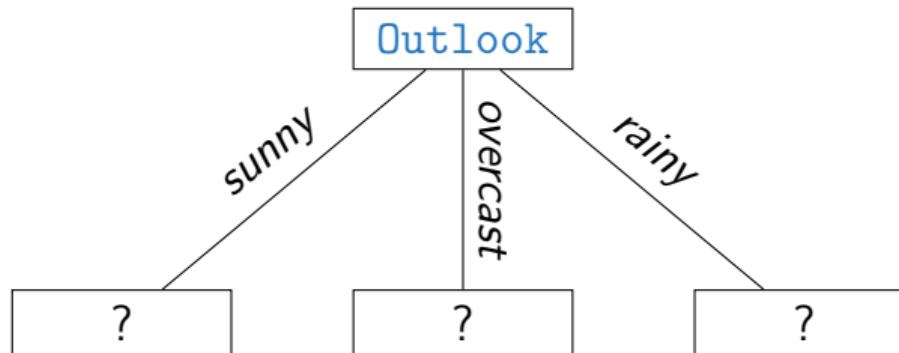
Outlook	Temperature	Humidity	Wind	PlayGolf
sunny	hot	high	weak	no
sunny	hot	high	strong	no
sunny	mild	high	weak	no
sunny	cool	normal	weak	yes
sunny	mild	normal	strong	yes
rainy	mild	high	weak	yes
rainy	cool	normal	weak	yes
rainy	cool	normal	strong	no
rainy	mild	normal	weak	yes
rainy	mild	high	strong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

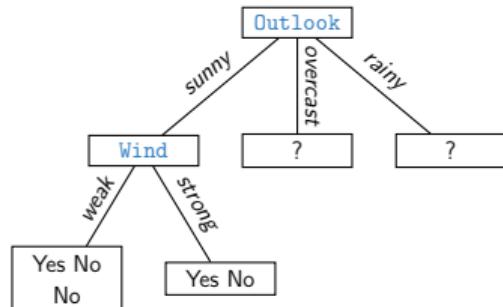
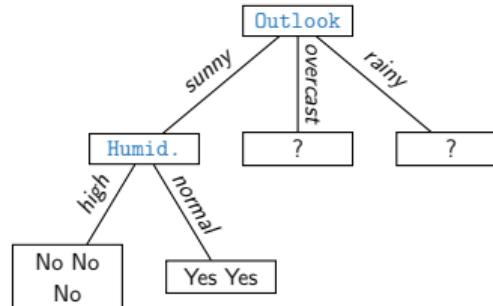
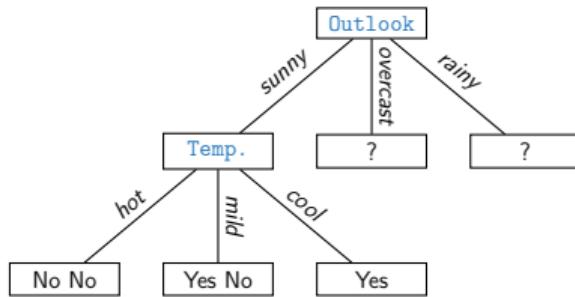
How to proceed?

- 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** \circlearrowleft (Go to 1)
- Recursion stops as soon as the subset is pure
- In the example above the subset $\mathcal{D}_{\text{outlook}=\text{overcast}}$ is already pure
- This algorithm is referred to as **ID3 (Iterative Dichotomizer)**

Step by Step: Construction of the Tree

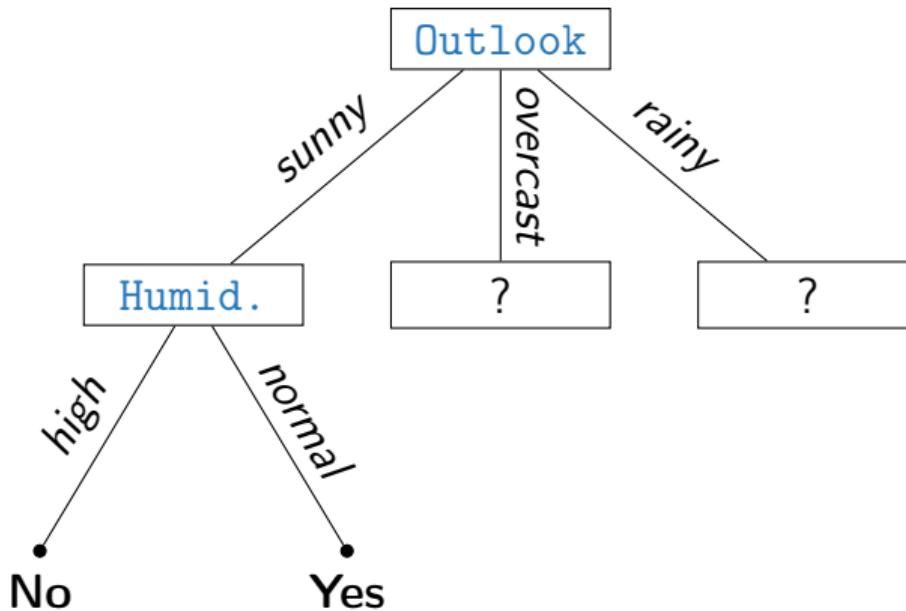


Step by Step: Construction of the Tree (Ctd.)

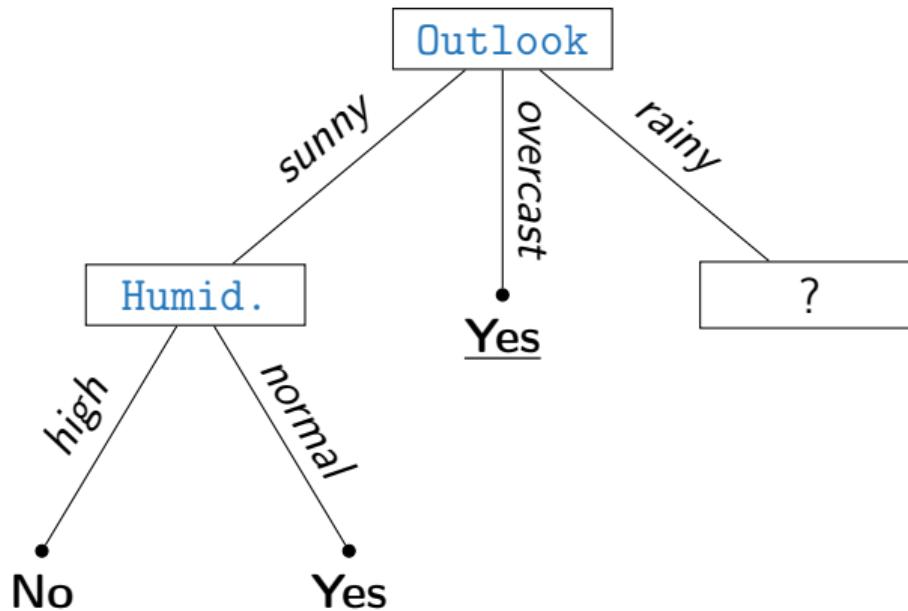


- $IG(\text{Temperature}) = 0.571$
- $IG(\text{Humidity}) = \mathbf{0.971}$
- $IG(\text{Wind}) = 0.020$

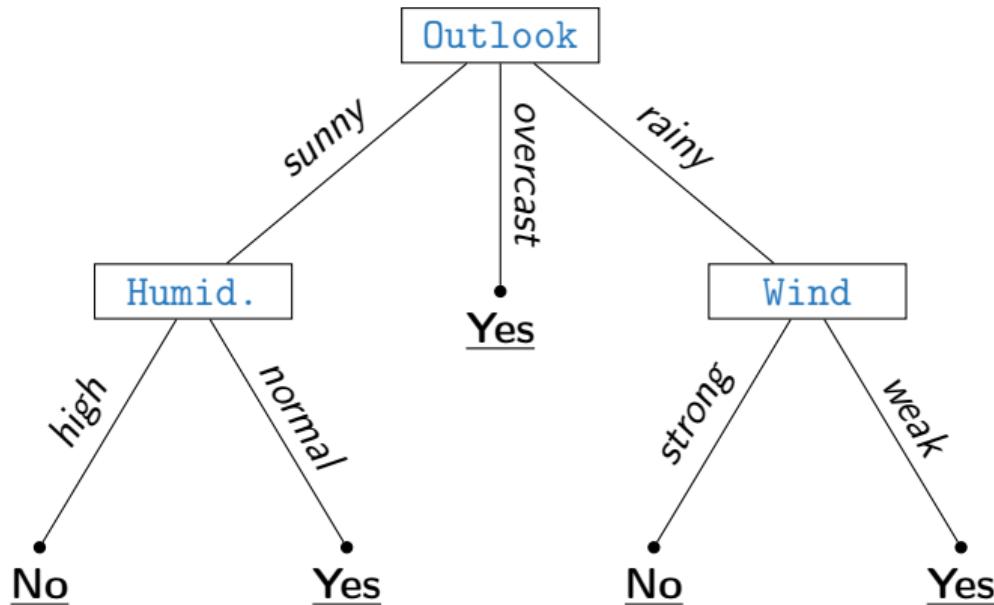
Step by Step: Construction of the Tree (Ctd.)



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



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
- 4 **if** $|Attr_List| = 0$ **then**
- 5 **return** N as leaf node labeled with majority class in \mathcal{D}
- 6 Find best split attribute A^* and label node N with A^*
- 7 $Attr_List \leftarrow Attr_List \setminus \{A^*\}$
- 8 **forall** $v \in dom(A^*)$ **do**
- 9 Let $\mathcal{D}_{A^*=v}$ be the set of tuples in \mathcal{D} that satisfy $A^* = v$
- 10 **if** $|\mathcal{D}_{A^*=v}| = 0$ **then**
- 11 Attach leaf labeled with majority class in \mathcal{D} to node N
- 12 **else**
- 13 Attach node returned by $ID3(\mathcal{D}_{A^*=v}, Attr_List)$
- 14 **return** N

Section:
Extensions and Variants



An Alternative to Information Gain: Gini Index

Gini index:

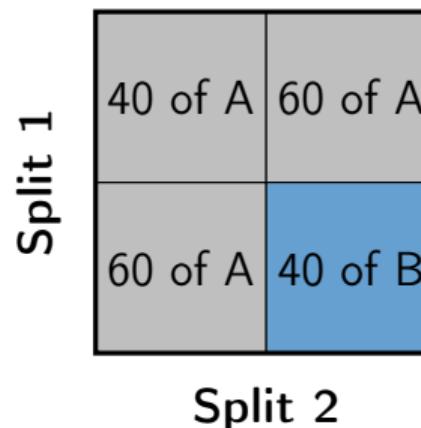
$$Gini(\mathcal{D}) = \sum_{c \in \mathcal{C}} p_c \cdot (1 - p_c) = 1 - \sum_{c \in \mathcal{C}} p_c^2 \quad (3)$$

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



Why not use the Error as a splitting Criterion?

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

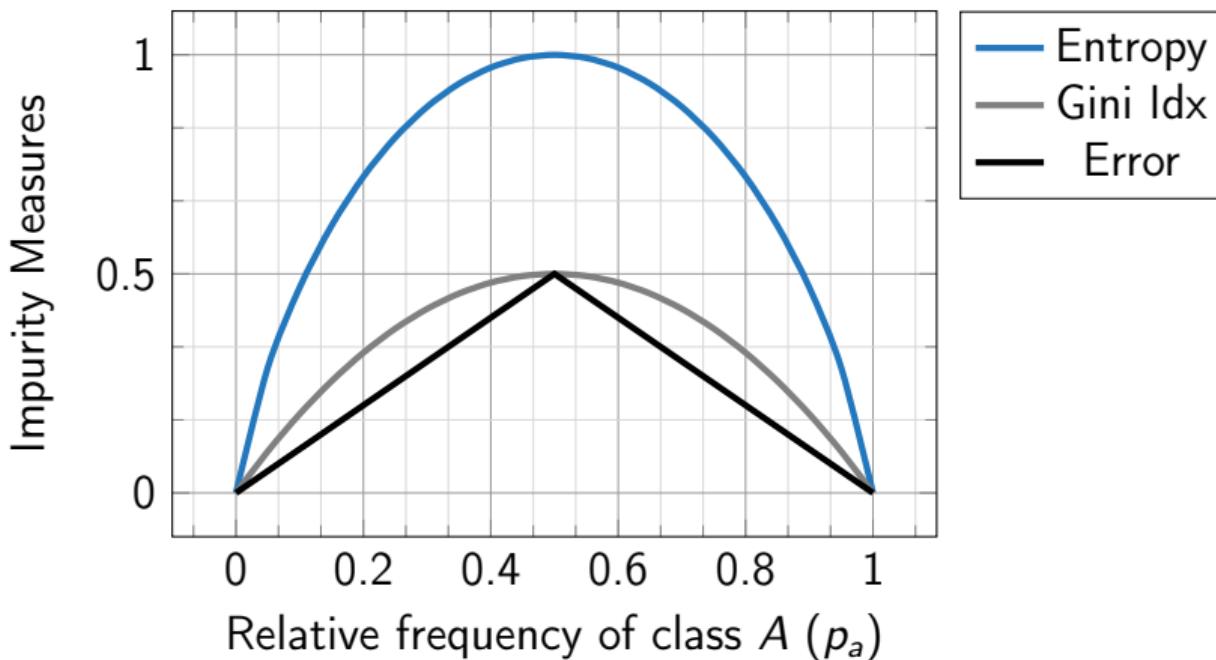


Error without splitting:
20 %

Error after splitting:
20 %

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

Summary: Impurity Measures



Highly-branching Attributes

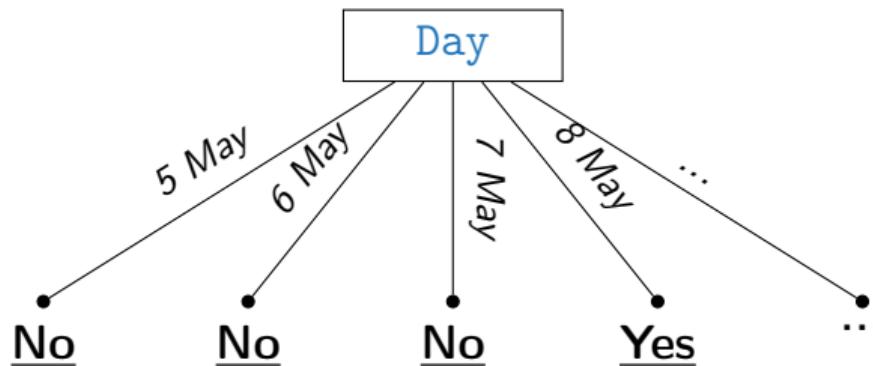
Attributes with a large number of values are problematic, since the leaves are not 'backed' with sufficient data examples.

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)

Highly-branching Attributes (Ctd.)



- Entropy before was 0.9403, Entropy after split is 0
- $IG(\mathcal{D}, \text{Day}) = 0.9403$
- Attribute Day would be chosen for the split \Rightarrow **Bad for prediction**

Highly-branching Attributes (Ctd.)

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

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

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

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

Highly-branching Attributes (Ctd.)

- Intrinsic information for attribute Day:

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

- Gain ratio for attribute Day:

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

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!**)

Handling numeric Attributes

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

Handling numeric Attributes (Ctd.)

- Consider the attribute Temperature:

Use **numerical values** instead of discrete values like *cool*, *mild*, *hot*:



- Temperature < 71.5

yes: 4 | no: 2

- Temperature ≥ 71.5

yes: 5 | no: 3

$$E(\mathcal{D}, \text{Temp.}) = \frac{6}{14} \cdot E(\text{Temp.} < 71.5) + \frac{8}{14} \cdot E(\text{Temp.} \geq 71.5) = \mathbf{0.939}$$

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
	\leq	>	\leq								
Yes	0	3	0	3	0	3	1	2	2	1	3
No	0	7	1	6	2	5	3	4	3	4	0
Gini	0.420	0.400	0.375	0.343	0.417	0.400	0.300	0.343	0.375	0.400	0.420



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):

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



Section:
Ensemble Methods



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
 - Stacking

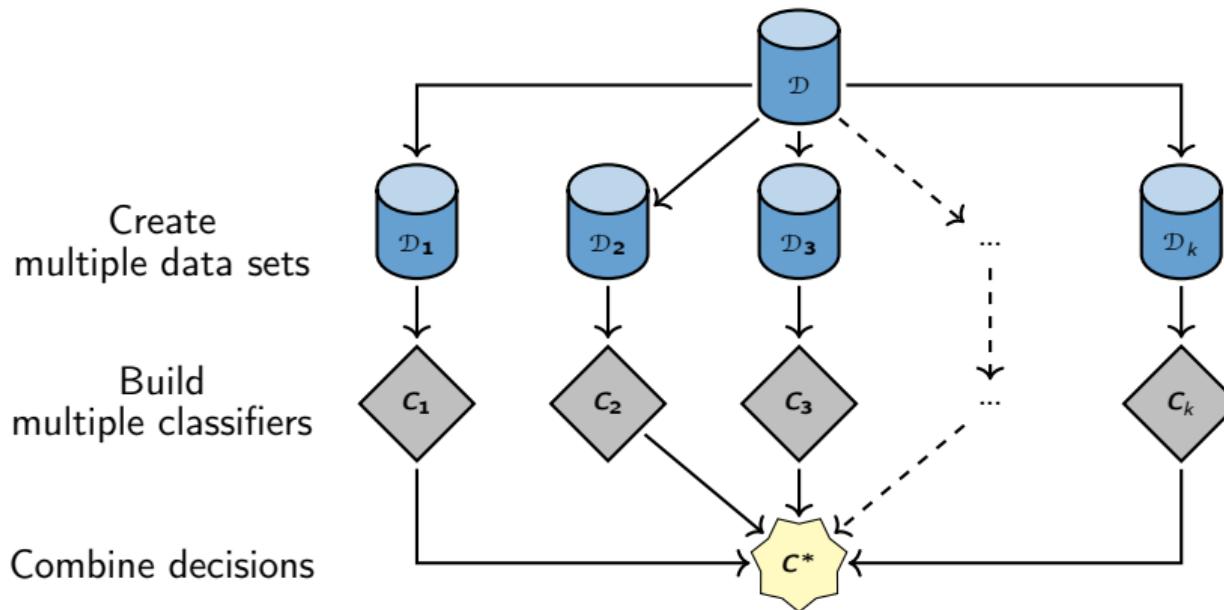
What is the Advantage?

- 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)

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

Bagging: General Approach

Bagging $\hat{=}$ Bootstrap Aggregating



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** $u \in \{1, 2, \dots, k\}$ **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^*)*

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** ⇒ **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, ...)

Bagged Decision Trees

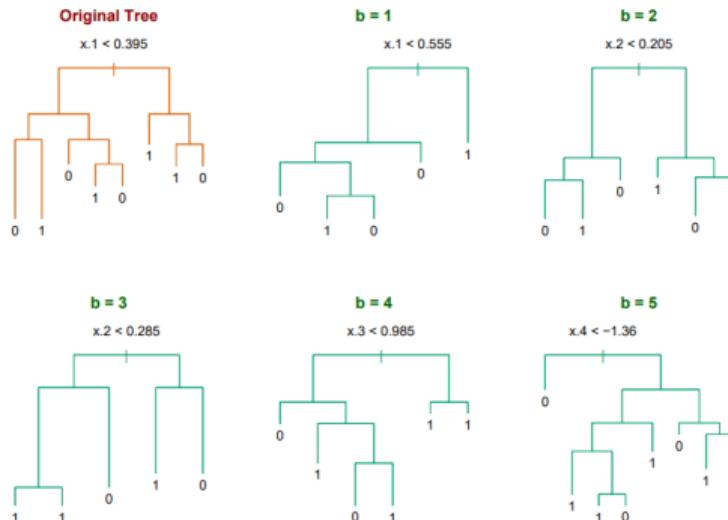


Figure: Bagged decision trees; cf. Hastie 2008, page 284

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.

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** $u \in \{1, 2, \dots, k\}$ **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^**)

ExtraTrees (Randomization 2.0)

- One more step of randomization ⇒ **Extremely Randomized 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.

Boosting Overview



- **Key Idea:**
 - New classifiers focus on examples misclassified by others
 - Assign a weight to each classifier (depending on their error)
- **How to:** (*unlike bagging, boosting **cannot** be parallelized!*)
 - ① Initialize example weights with $1/n$
 - ② Train a model using the example weights, calculate error and model weight
 - ③ Update example weights according to model weight
 - ④ Train the next model using the updated example weights
 - ⑤ Combine predictions according to model weights

Algorithm 4: AdaBoost Algorithm

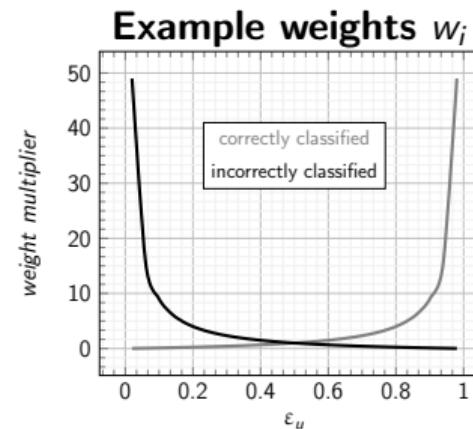
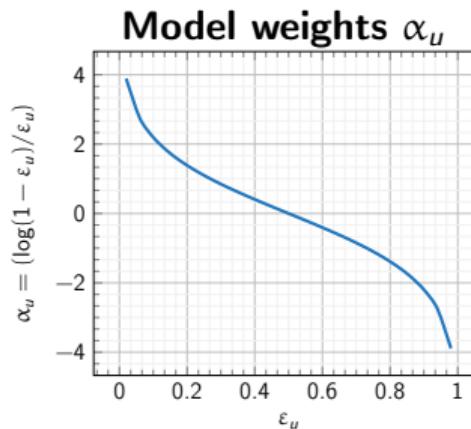
Input: Training set \mathcal{D} , number of base classifiers k

- 1 **Training:**
 - 2 Initialize example weights $w_i \leftarrow \frac{1}{n}$
 - 3 **forall** $u \in \{1, 2, \dots, k\}$ **do**
 - 4 Learn a classifier C_u using the current example weights
 - 5 Compute weighted error estimate of the model: $\varepsilon_u \leftarrow \sum_{i=1}^n w_i \cdot \mathbb{1}\{y_{pred} \neq y^{(i)}\}$
 - 6 Compute model weight: $\alpha_u \leftarrow 1/2 \cdot \ln\left(\frac{1-\varepsilon_u}{\varepsilon_u}\right)$
 - 7 For all correctly classified examples: $w_i \leftarrow w_i \cdot \exp\{-\alpha_u\}$
 - 8 For all incorrectly classified examples: $w_i \leftarrow w_i \cdot \exp\{\alpha_u\}$
 - 9 Normalize the weights w_i (such that they sum up to one): $\frac{w_i}{\sum_{i=1}^n w_i}$

 - 10 **Prediction:**
 - 11 **forall** *unlabeled instances* **do**
 - 12 Get predictions from all classifiers C_u
 - 13 **return** Class which receives the highest sum of weights α_u (combined classifier C^*)
-



AdaBoost: Illustration of Weights



- Classifier weights α_u grow exponentially
- Classifier weight is 0 if error is $1/2$ (model cannot be trusted, random guessing)
- For high errors \Rightarrow **Do the opposite of what the classifiers predicts**

AdaBoost: Two Classes

- With two classes, -1 and $+1$, the final classifier C^* can be written as:

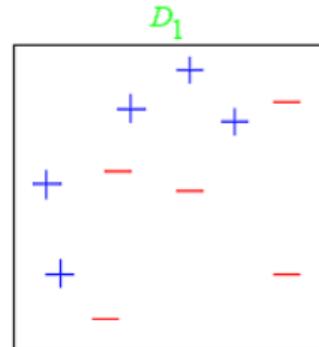
$$H(\mathbf{x}) = \text{sign} \left(\sum_{u=1}^k \alpha_u h_u(\mathbf{x}) \right) \quad (10)$$

- α_u is defined as above
- The weight update for the examples is given by:

$$w_i^{new} = \frac{w_i^{old} \cdot \exp\{-\alpha_u y^{(i)} h_i(\mathbf{x}^{(i)})\}}{\sum_{i=1}^n w_i} \quad (11)$$

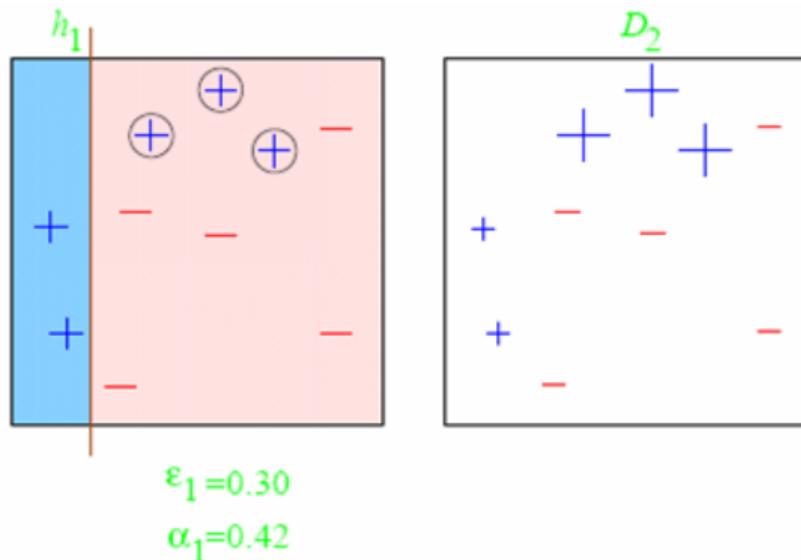
Boosting Example

- Ensemble of **decision stumps** ($\hat{=}$ tree with a single split only)
- 10 examples (5 from class +1 and 5 from class -1)
- Each example is initially weighted with a factor of $1/10$



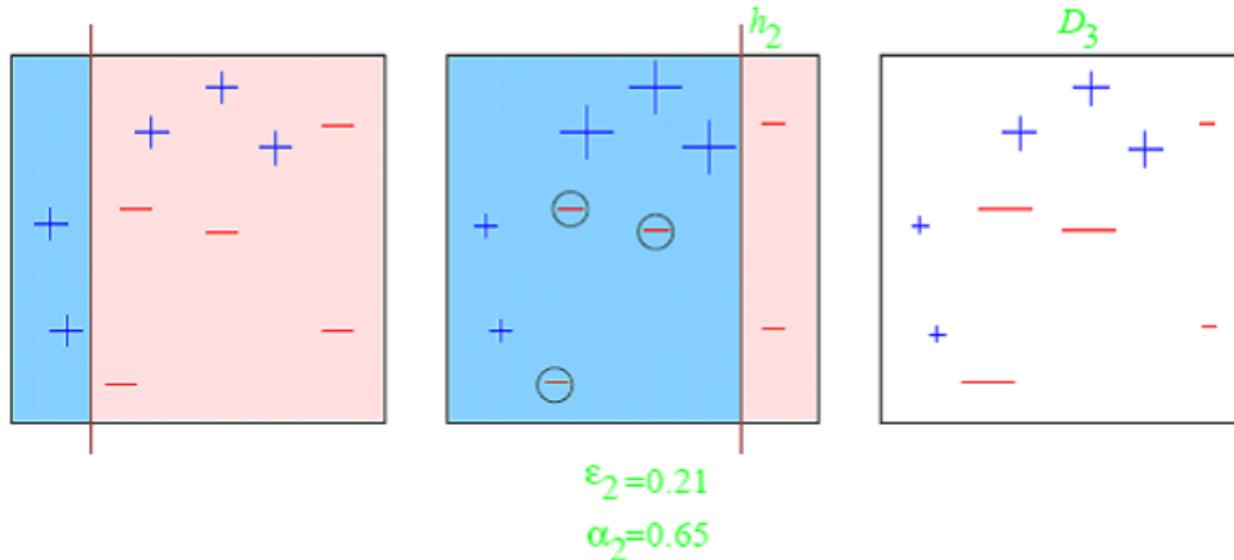
Boosting Example (Ctd.)

1st Iteration



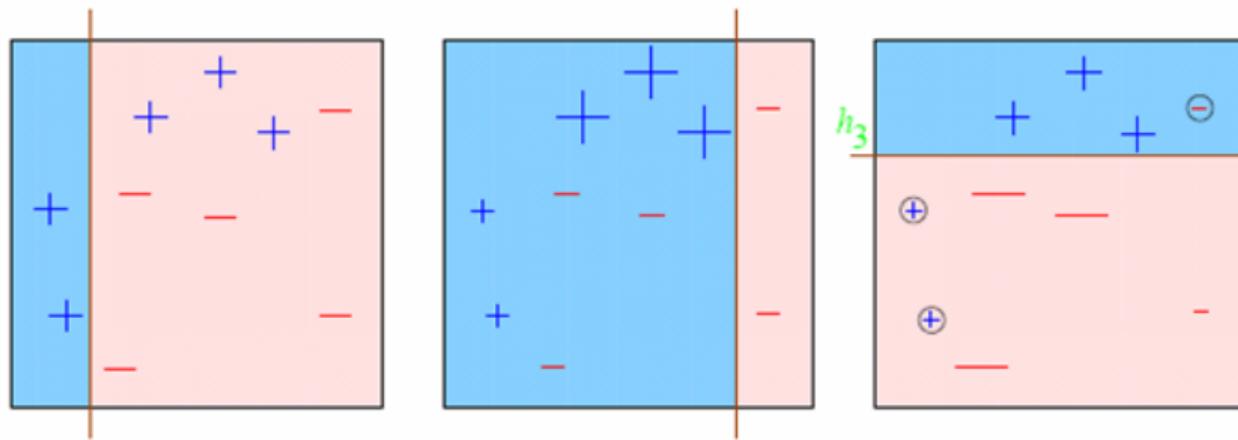
Boosting Example (Ctd.)

2nd Iteration



Boosting Example (Ctd.)

3rd Iteration

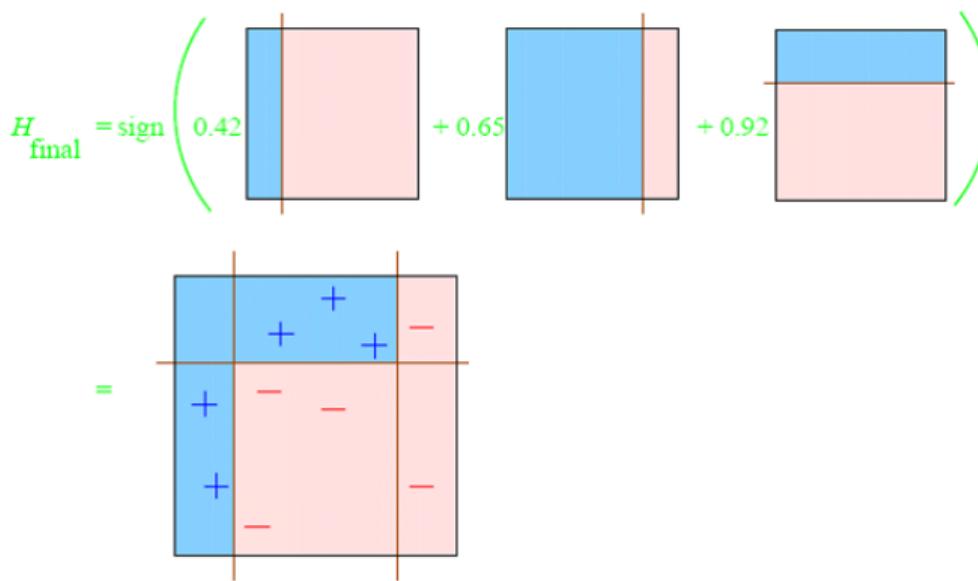


$$\varepsilon_3 = 0.14$$

$$\alpha_3 = 0.92$$

Boosting Example (Ctd.)

Prediction



Combining Predictions

① Voting ✓

- Each classifier votes for one of the classes
- Each classifier has the same (model) weight
- This is **Bagging**

② Weighted Voting ✓

- Weight the individual predictions with model weight
- This is **Boosting**

③ Stacking

- Use a **meta classifier** to combine the predictions
- See next slides...

Algorithm 5: Stacking Algorithm

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

- 1 Train k base classifiers C_1, C_2, \dots, C_k
 - 2 Let all base classifiers predict the labels for the training set
 - 3 Create the meta data set (*)
 - 4 Train a meta classifier M on the meta data set
-

(*) Meta data set:

- **Labels:**

Use the same labels as in the original data set

- **Attributes:**

One attribute for each base classifier (k attributes)

The attribute values are the predictions of the corresponding classifiers

Stacking Example

Attributes		Class
$x_1^{(1)}$...	$x_m^{(1)}$
$x_1^{(2)}$...	$x_m^{(2)}$
...
$x_1^{(n)}$...	$x_m^{(n)}$
		<i>true</i>
		<i>false</i>
		...
		<i>true</i>

Table: Original training set

C_1	C_2	...	C_k
<i>true</i>	<i>true</i>	...	<i>false</i>
<i>false</i>	<i>true</i>	...	<i>true</i>
...
<i>false</i>	<i>false</i>	...	<i>true</i>

Table: Predictions base classifiers

C_1	C_2	...	C_k	Class
<i>true</i>	<i>true</i>	...	<i>false</i>	<i>true</i>
<i>false</i>	<i>true</i>	...	<i>true</i>	<i>false</i>
...
<i>false</i>	<i>false</i>	...	<i>true</i>	<i>true</i>

Table: Meta data set



Multi-Class-Problems

- Usually, a **class binarization** technique is needed
- One-vs-Rest (one classifier for each class)
- Each class is encoded in a **bit-string**

Class	Class Vector			
a	1	0	0	0
b	0	1	0	0
c	0	0	1	0
d	0	0	0	1

What would you do if the result were **1100**? It is not clear if the new instance belongs to class **a** or class **b**.
Solution: Use error-correcting output codes instead!



Error-Correcting Output Codes (ECOCs)

- Use code words that have a high pairwise **Hamming distance d**
- Can correct up to $(d - 1)/2$ single bit-errors
- In this case: Use seven classifiers

Class	Class Vector						
a	1	1	1	1	1	1	1
b	0	0	0	0	1	1	1
c	0	0	1	1	0	0	1
d	0	1	0	1	0	1	0

- Predicted code-word: **1011111** ⇒ What is the true class???
Probably class a, the second classifier made a mistake...



ECOCs: How to encode the classes?



- Criteria for the code words
 - **Row separation** guarantees that errors can be corrected
 - **Column separation** (identical columns \Rightarrow Classifier makes same error several times)
- **Exhaustive Codes**
 - **Class 1:** Code word consists of 1s only
 - **Class 2:** 2^{k-2} 0s followed by $2^{k-2} - 1$ 1s ($k \hat{=} \#$ classes)
 - **Class i:** Alternating runs of 2^{k-i} 0s and 1s
(Last run is one bit shorter than the others)

Section:
Wrap-Up



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:**
 - **Bagging:** Sample diverse data sets from underlying data
 - **Boosting:** Focus on instances misclassified by earlier classifiers
 - **Stacking:** Learn a meta-classifier on top of base classifiers
 - **Random forests** combine bagging with randomization

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 distinguishes bagging from boosting?
- ⑦ Explain what a random forest does.

What's next...?

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

Recommended Literature and further Reading I



[1] The Elements of Statistical Learning

Hastie et al. Springer. 2008.

→ [Link](#), cf. chapters 8.7, 8.8, 9.2, 10.1



[2] Machine Learning

Tom M. Mitchell. McGraw-Hill. 1997.

→ [Link](#), cf. chapter 3



[3] Using Bagging and Boosting to Improve Classification Tree Accuracy

Dave Sotelo, Online Blog

→ [Link](#)

Recommended Literature and further Reading II



[4] Pattern Recognition and Machine Learning

Christopher Bishop. Springer. 2006.

→ [Link](#), cf. chapter 14.3

Meme of the Day



Thank you very much for the attention!

Topic: *** Applied Machine Learning Fundamentals *** Decision Trees and Ensembles

Term: Winter term 2019/2020

Contact:

M. Sc. Daniel Wehner

SAP SE

daniel.wehner@sap.com

Do you have any questions?