

Lecture 2: Decision Trees

Andreas Wichert
Department of Computer Science and Engineering
Técnico Lisboa

1

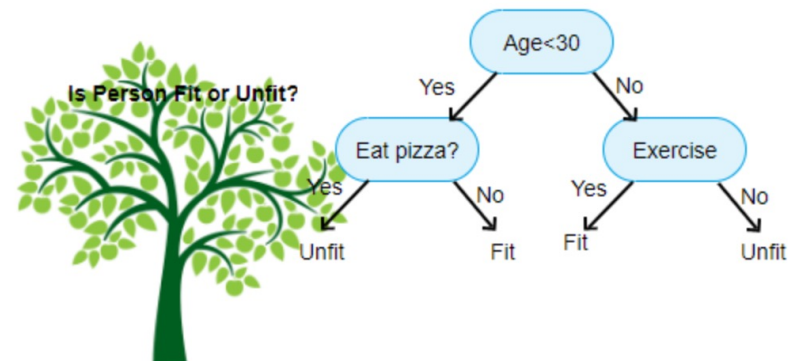
- Learning through search (AI)
 - Blind search
 - Impossible since the search space grows exponentially fast
 - Greedy Search
 - Doesn't guarantee to find the best result
- Relation to Gradient Descent, local minima problem....

2

When to consider Decision Trees

- Instances describable by attribute-value pairs
- Target function is discrete valued
- Possibly noisy training data
- Missing attribute values
- Examples:
 - Data Bases/ Data Warehouse Analysis
 - Credit risk analysis
 - Medical diagnosis
 - Object classification for robot manipulator (Tan 1993)

3



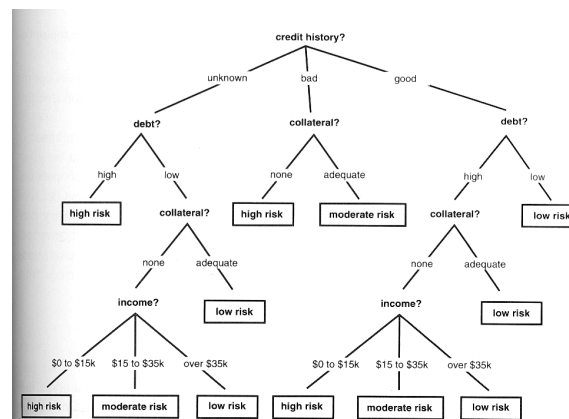
4

TABLE 12.1 DATA FROM CREDIT HISTORY OF LOAN APPLICATIONS

NO.	RISK	CREDIT HISTORY	DEBT	COLLATERAL	INCOME
1.	high	bad	high	none	\$0 to \$15k
2.	high	unknown	high	none	\$15 to \$35k
3.	moderate	unknown	low	none	\$15 to \$35k
4.	high	unknown	low	none	\$0 to \$15k
5.	low	unknown	low	none	over \$35k
6.	low	unknown	low	adequate	over \$35k
7.	high	bad	low	none	\$0 to \$15k
8.	moderate	bad	low	adequate	over \$35k
9.	low	good	low	none	over \$35k
10.	low	good	high	adequate	over \$35k
11.	high	good	high	none	\$0 to \$15k
12.	moderate	good	high	none	\$15 to \$35k
13.	low	good	high	none	over \$35k
14.	high	bad	high	none	\$15 to \$35k

5

Decision tree for credit risk assessment



6

- The decision tree represents the classification of the table
- It can classify all the objects in the table
- Each internal node represents a test on some property
- Each possible value of that property corresponds to a branch of the tree
- An individual of unknown type may be classified by traversing this tree

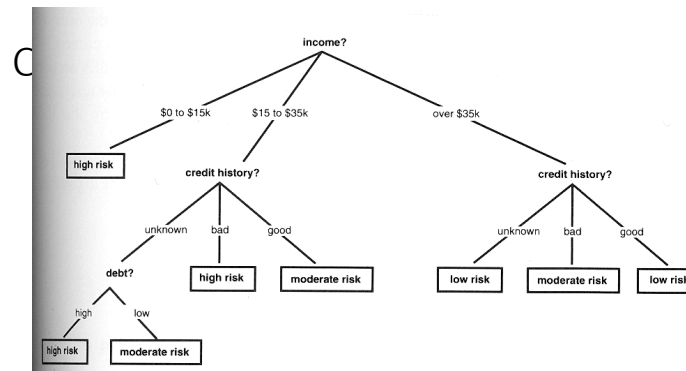
7

- In classifying any given instance, the tree **does not use all the properties** in the table
- Decision tree for credit risk assessment
- If a person has a good credit history and low debit, we ignore her collateral income and classify her as low risk
- In spite of omitting certain tests, the tree classifies all examples in the table

8

- In general, the size of a tree necessary to classify a given set of examples **varies according to the order** with which properties (=attributes) are tested
- Given a set of training instances and a number of **different** decision trees that correctly classify the instances, we may ask which tree has the greatest likelihood of correctly classifying using instances of the population?

9

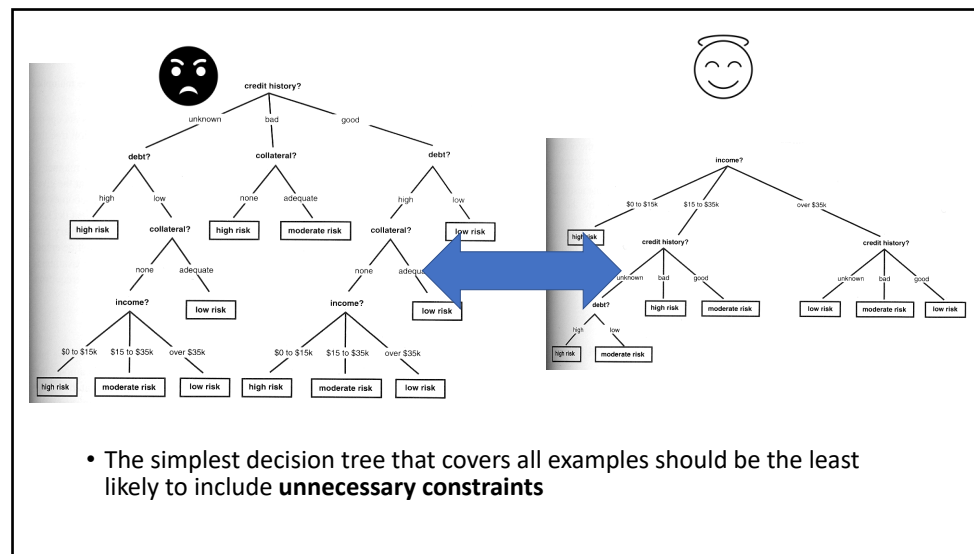


- This is a simplified decision tree for credit risk assessment
 - It classifies all examples of the table correctly

10

- ID3 algorithm assumes that a **good decision tree is the simplest decision tree**
- Heuristic:
 - Preferring simplicity **and avoiding unnecessary assumptions**
 - Known as Occam's Razor
- The simplest decision tree that covers all examples should be the least likely to include **unnecessary constraints**

11



12

Occam Razor



- Occam Razor was first articulated by the medieval logician William of Occam in 1324
 - *born in the village of Ockham in Surrey (England) about 1285, believed that he died in a convent in Munich in 1349, a victim of the Black Death*
 - *It is vain do with more what can be done with less..*
- We should always accept the **simplest answer** that correctly fits our data
- The **smallest decision tree** that correctly classifies all given examples

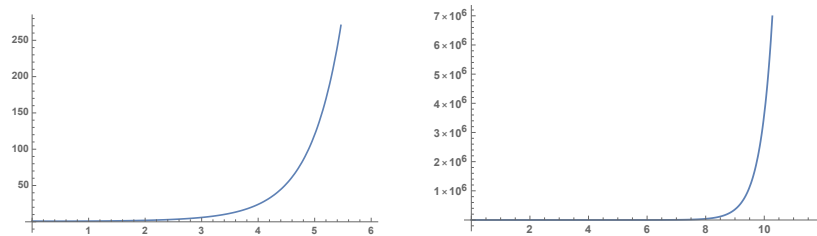
13

- Because the order of tests is critical to constructing a simple tree, ID3 relies heavily on its criteria for selecting the test at the root of each sub tree
- How many different decision tree exist?
- N =Number Of Attributes
- There exist $N!$ different ordering of attributes, different decision trees
- Algorithm:
 - Blind Search finds the global minima, the smallest decision tree (optimal)
 - Compute all $N!$ decision trees and chose the smallest one

14

$N!$ grows extremely fast

- $N=4$, $4!=24$ manageable, no problem



- Plot of $N!$

15

We have to use a **heuristic function**

- ID3 selects a property **heuristic** to test at the current node of the tree and uses this test to partition the set of examples
- The algorithm then recursively constructs a sub tree for each partition
- This continues until all members of the partition are in the same class
- That class becomes a leaf node of the tree

16

Top-Down Induction of Decision Trees ID3

1. $A \leftarrow$ the "best" decision attribute for next *node*
2. Assign A as decision attribute (=property) for *node*
3. For each value of A create new descendant
4. Sort training examples to leaf node according to the attribute value of the branch
5. If all training examples are perfectly classified (same value of target attribute) stop, else iterate over new leaf nodes

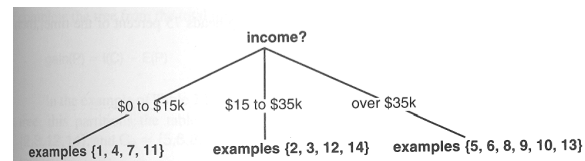
17

TABLE 12.1 DATA FROM CREDIT HISTORY OF LOAN APPLICATIONS

NO.	RISK	CREDIT HISTORY	DEBT	COLLATERAL	INCOME
1.	high	bad	high	none	\$0 to \$15k
2.	high	unknown	high	none	\$15 to \$35k
3.	moderate	unknown	low	none	\$15 to \$35k
4.	high	unknown	low	none	\$0 to \$15k
5.	low	unknown	low	none	over \$35k
6.	low	unknown	low	adequate	over \$35k
7.	high	bad	low	none	\$0 to \$15k
8.	moderate	bad	low	adequate	over \$35k
9.	low	good	low	none	over \$35k
10.	low	good	high	adequate	over \$35k
11.	high	good	high	none	\$0 to \$15k
12.	moderate	good	high	none	\$15 to \$35k
13.	low	good	high	none	over \$35k
14.	high	bad	high	none	\$15 to \$35k

18

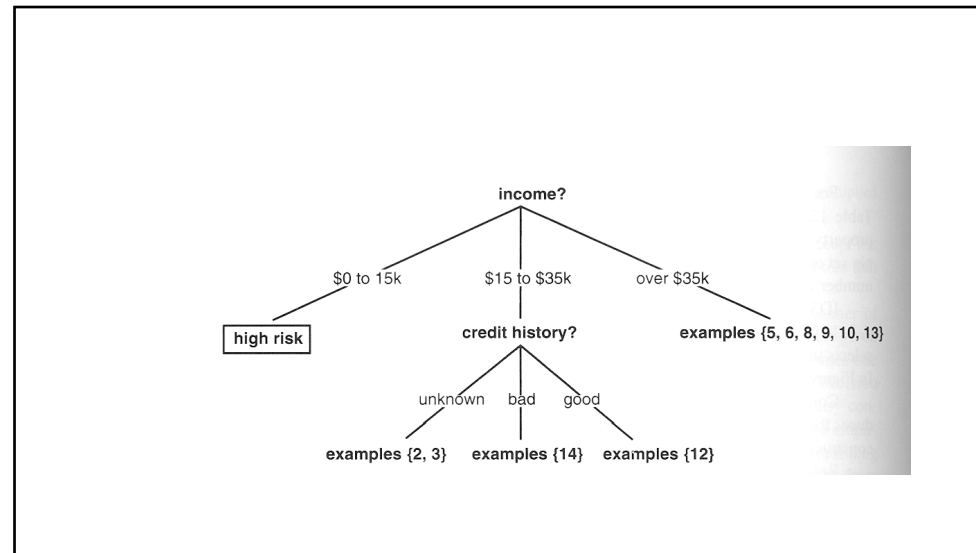
- ID3 constructs the tree for credit risk assessment
 - Beginning with the full table of examples, ID3 selects **income** as the root property using function selecting “best” property (attribute)
 - The examples are divided, listed by their number in the list



19

- ID3 applies the method recursively for each partition
 - The partition {1,4,7,11} consists entirely of high-risk individuals, a node is created
 - ID3 selects credit history property as the root of the subtree for the partition {2,3,12,14}
 - Credit history further divides this partition into {2,3},{14} and {12}
- ID3 implements a form of hill climbing in the space of all possible trees using a heuristic function
- Doesn't guarantee to find the smallest decision tree, can find a local maxima.

20



21

Heuristic function: Shannon Entropy

- Shannon formalized these intuitions
- Given a universe of messages $M=\{m_1, m_2, \dots, m_n\}$ and a probability $p(m_i)$ for the occurrence of each message, the **information** content (also called entropy) of a message M is given

$$I(M) = \sum_{i=1}^n -p(m_i) \log_2(p(m_i))$$

22



- Claude Elwood Shannon (1916 - 2001) was an American mathematician, electrical engineer, and cryptographer known as “the father of information theory”

23

$$H = - \sum_t^K p(x_t) \cdot \log_2 p(x_t) = - \sum_{x \in X} p(x) \cdot \log_2 p(x).$$

N different events x_t and $p(x_t)$ as the probability of occurrence of the events, the theoretical minimum average number of bits is computed using Shannon's formula of entropy

24

- Information content of a message telling the outcome of the flip of an honest coin

$$I(\text{Coin_toss}) = -p(\text{heads})\log_2(p(\text{heads})) - p(\text{tails})\log_2(p(\text{tails}))$$

$$I(\text{Coin_toss}) = -p(0.5)\log_2(p(0.5)) - p(0.5)\log_2(p(0.5))$$

$$I(\text{Coin_toss}) = 1 \quad \text{bit}$$

25

- However if the coin has been rigged to come up heads 75 percent

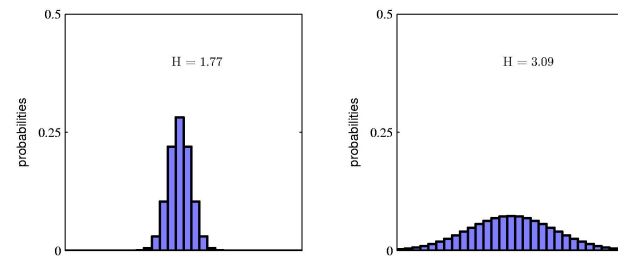
$$I(\text{Coin_toss}) = -p(\text{heads})\log_2(p(\text{heads})) - p(\text{tails})\log_2(p(\text{tails}))$$

$$I(\text{Coin_toss}) = -p(0.75)\log_2(p(0.75)) - p(0.25)\log_2(p(0.25))$$

$$I(\text{Coin_toss}) = 0.811 \quad \text{bits}$$

26

Shannon Entropy



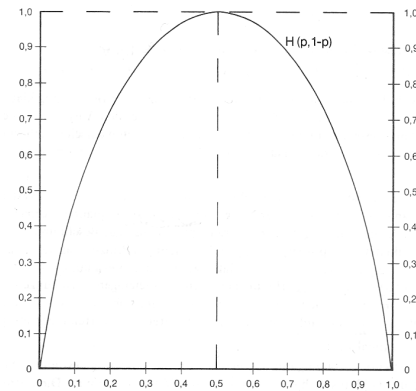
27

- An experiment is described by probabilities $p=(p_1, p_2, \dots, p_n)$
- Does the distribution of these probabilities have an effect on the entropy?
- It turns out that the entropy is maximal in the case all probabilities are equal, means $p=(1/n, 1/n, \dots, 1/n)$
- In this case the maximal ideal Entropy is

$$H(F) = - \sum_i p_i \log_2 p_i = - \log_2 1/n = \log_2 n$$

28

Only two probabilities



29

- The relationship between \log_2 and any other base b involves multiplication by a constant,

$$\log_2 x = \frac{\log_b x}{\log_b 2} = \frac{\log_{10} x}{\log_{10} 2}.$$

$$H = -\frac{1}{\log_{10} 2} \cdot \sum_i^n p(m_i) \cdot \log_{10} p(m_i) = -\sum_i^n p(m_i) \cdot \log_2 p(m_i)$$

30

- We may think of a decision tree as conveying information about the classification of examples in the decision table
- The information content of the tree is computed from the *probabilities* of different classifications.
 - If Ω is the set of all possible events, $p(\Omega) = 1$, then $a \in \Omega$.
 - $card(\Omega)$ is the number of elements of the set Ω , $card(a)$ is the number of elements of the set a and

$$p(a) = \frac{card(a)}{card(\Omega)}$$

31

TABLE 12.1 DATA FROM CREDIT HISTORY OF LOAN APPLICATIONS

NO.	RISK	CREDIT HISTORY	DEBT	COLLATERAL	INCOME
1.	high	bad	high	none	\$0 to \$15k
2.	high	unknown	high	none	\$15 to \$35k
3.	moderate	unknown	low	none	\$15 to \$35k
4.	high	unknown	low	none	\$0 to \$15k
5.	low	unknown	low	none	over \$35k
6.	low	unknown	low	adequate	over \$35k
7.	high	bad	low	none	\$0 to \$15k
8.	moderate	bad	low	adequate	over \$35k
9.	low	good	low	none	over \$35k
10.	low	good	high	adequate	over \$35k
11.	high	good	high	none	\$0 to \$15k
12.	moderate	good	high	none	\$15 to \$35k
13.	low	good	high	none	over \$35k
14.	high	bad	high	none	\$15 to \$35k

32

- The credit history loan table has following information (with $\Omega=14$)

- $p(\text{risk is high})=6/14$
- $p(\text{risk is moderate})=3/14$
- $p(\text{risk is low})=5/14$

$$p(a) = \frac{\text{card}(a)}{\text{card}(\Omega)}$$

$$I(\text{credit_table}) = -\frac{6}{14} \log_2\left(\frac{6}{14}\right) - \frac{3}{14} \log_2\left(\frac{3}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right)$$

$$I(\text{credit_table}) = 1.531 \quad \text{bits}$$

33

- For a given test, the information gain provided by making that test at the root of the current tree is equal to
- Total information of the table - the amount of information needed to complete the classification after performing the test
- The amount of information needed to complete the tree is defined as weighted average of the information content of each sub tree

34

- The amount of information needed to complete the tree is defined as **weighted average** of the information content of each sub tree by the percentage of the examples present
- C a set of training instances. If property (for example *income*) with n values, C will be divided into the **subsets** $\{C_1, C_2, \dots, C_n\}$
- Expected information needed to complete the tree after making P root

$$E(P) = \sum_{i=1}^n \frac{|C_i|}{|C|} I(C_i)$$

35

- In the credit history loan table we make **income** the property tested at the root
- This makes the division into

- $C_1 = \{1, 4, 7, 11\}$
 - $I(C_1) = -4/4 * \log_2(4/4) = 0$ bit
- $C_2 = \{2, 3, 12, 14\}$
 - $I(C_2) = -1/2 * \log_2(1/2) - 1/2 * \log_2(1/2) = 1$ bit
- $C_3 = \{5, 6, 8, 9, 10, 13\}$
 - $I(C_3) = -5/6 * \log_2(5/6) - 1/6 * \log_2(1/6) = 0.650022$ bit

$$E(\text{income}) = \frac{4}{14} I(C_1) + \frac{4}{14} I(C_2) + \frac{6}{14} I(C_3)$$

$$E(\text{income}) = \frac{4}{14} 0 + \frac{4}{14} 1.0 + \frac{6}{14} 0.65$$

$$E(\text{income}) = 0.564 \quad \text{bits}$$

TABLE 12.1 DATA FROM CREDIT HISTORY OF LOAN APPLICATIONS

NO.	RISK	CREDIT HISTORY	DEBT	COLLATERAL	INCOME
1.	high	bad	high	none	\$0 to \$15k
2.	high	unknown	high	none	\$15 to \$35k
3.	moderate	unknown	low	none	\$15 to \$35k
4.	high	unknown	low	none	\$0 to \$15k
5.	low	unknown	low	none	over \$35k
6.	low	unknown	low	adequate	over \$35k
7.	high	bad	low	none	\$0 to \$15k
8.	moderate	bad	low	adequate	over \$35k
9.	low	good	low	none	over \$35k
10.	low	good	high	adequate	over \$35k
11.	high	good	high	none	\$0 to \$15k
12.	moderate	good	high	none	\$15 to \$35k
13.	low	good	high	none	over \$35k
14.	high	bad	high	none	\$15 to \$35k

36

$$E(P) = \sum_{i=1}^n \frac{|C_i|}{|C|} I(C_i)$$

- The gain from the property P is computed by subtracting the expected information to complete E(P) from the total information

$$\text{gain}(P) = I(C) - E(P)$$

37

$$\text{gain}(\text{income}) = I(\text{credit_table}) - E(\text{income})$$

$$\text{gain}(\text{income}) = 1.531 - 0.564$$

$$\text{gain}(\text{income}) = \mathbf{0.967 \text{ bits}}$$

$$\text{gain}(\text{credit history}) = 0.266$$

$$\text{gain}(\text{debt}) = 0.581$$

$$\text{gain}(\text{collateral}) = 0.756$$

38

- Because income provides the greatest information gain, ID will select it as the root of the tree
- The algorithm continues to apply this analysis recursively to each **subtree**, until it has completed the tree.

39

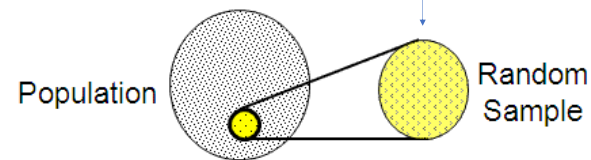
Overfitting

- The ID3 algorithm grows each branch of the tree just deeply enough to perfectly classify the training examples
- Difficulties may be present:
 - When there is noise in the data
 - When the number of training examples is too small to produce a representative sample of the true target function
- The ID3 algorithm can produce trees that **overfit** the training examples

40

Overfitting

- We will say that a hypothesis overfits the training examples - if some other hypothesis that fits the training examples *less well* actually performs better over the *entire* distribution of instances (included instances beyond training set)



41

Overfitting

Consider error of hypothesis h over

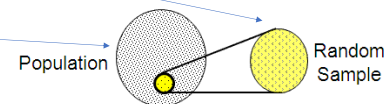
- Training data: $error_{train}(h)$
- Entire distribution D of data: $error_D(h)$

Hypothesis $h \in H$ *overfits* training data if there is an alternative hypothesis $h' \in H$ such that

$$error_{train}(h) < error_{train}(h')$$

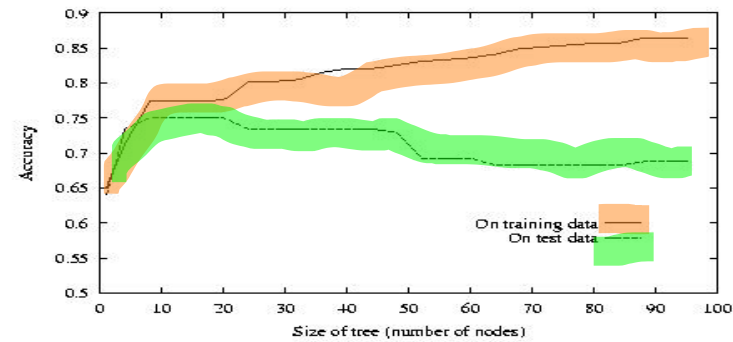
and

$$error_D(h) > error_D(h')$$



42

Overfitting



43

- How can it be possible for a tree h to fit the *training examples* better than h' , but to perform more poorly over *subsequent examples*
- One way this can occur when *the training examples* contain random errors or noise

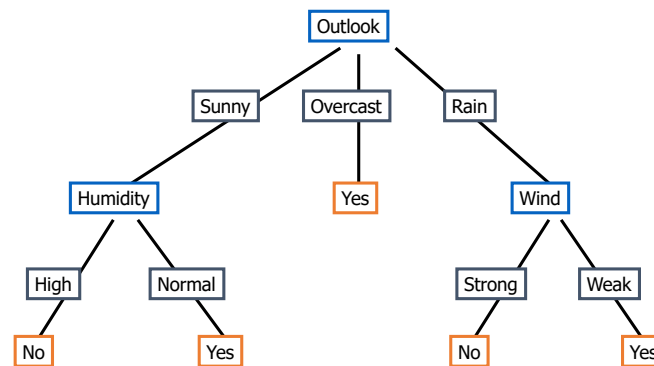
44

Training Examples

Day	Outlook	Temp.	Humidity	Wind	Play Tennis
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	Weak	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cold	Normal	Weak	Yes
D10	Rain	Mild	Normal	Strong	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

45

Decision Tree for PlayTennis



46

- Consider of adding the following positive training example, **incorrectly** labeled as negative
- *Outlook=Sunny, Temperature=Hot, Humidly=Normal, Wind=Strong, PlayTennis=No*
- The addition of this **incorrect example** will now cause ID3 to construct a more complex tree
- Because the new example is labeled as a negative example, ID3 will search for further refinements to the tree

47

- As long as the new erroneous example differs in some attributes, ID3 will succeed in finding a tree
- ID3 will output a decision tree (h) that is **more complex** than the original tree (h')
- Given the new decision tree a simple consequence of fitting noisy training examples, h' will outperform h on the **test set**

48

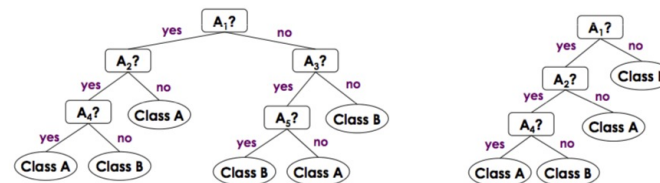
Avoid Overfitting

- How can we avoid overfitting?
 - Stop **growing** when data split not statistically significant
 - Grow full tree then post-prune
- How to select ``best" tree:
 - Measure performance over training data
 - Measure performance over separate validation data set

49

Pruning

- Remove the least reliable branches



50

Committees

- The simplest way to construct a committee is to **average the predictions** of a set of individual models
- We have only a single data set, and so we have to find a way to introduce variability between the different models within the committee.

51

Bootstrap a data set

- Sampling a dataset with replacement
- Define: Size of the sample and the number of repeats.
- Example:
 - $(0.1, 0.2, 0.3, 0.4, 0.5, 0.6)$
 - Randomly choose the first observation from the dataset
 - $sample = (0.2)$
- This observation is returned to the dataset and we repeat this step 3 more times.
 - $sample = (0.2, 0.1, 0.2, 0.6)$

52

Bagging

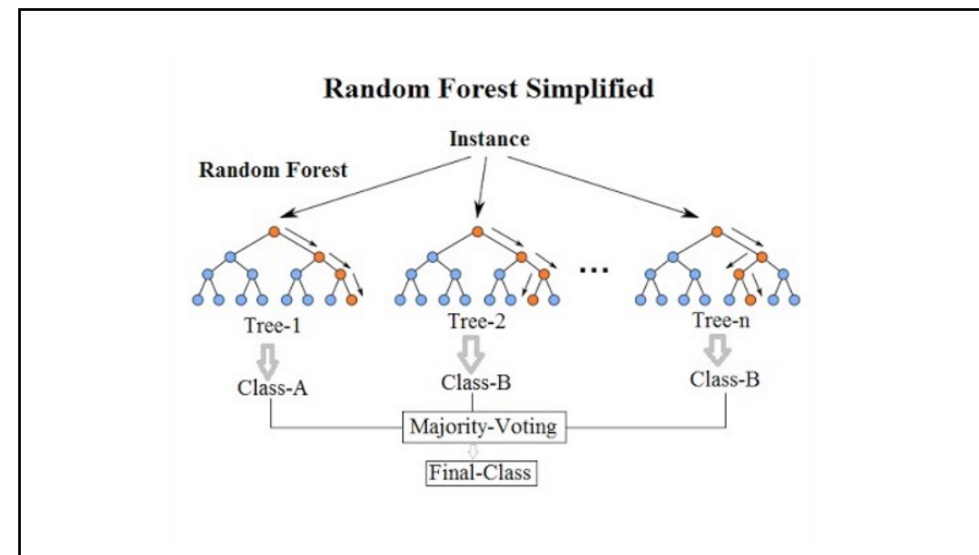
Suppose we generate M bootstrap data sets and then use each to train a separate copy $y_m(\mathbf{x})$ of a predictive model where $m = 1, \dots, M$.

The committee prediction is given by

$$y_{COM}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M y_m(\mathbf{x})$$

This procedure is known as bootstrap aggregation or bagging

53



54

Quinlan strategies of C4.5

- Derive an initial rule set by enumerating paths from the root to the leaves
- Generalize the rules by possible deleting conditions deemed to be unnecessary
- Group the rules into subsets according to the target classes they cover
 - Delete any rules that do not appear to contribute to overall performance on that class
- Order the set of rules for the target classes, and chose a default class to which cases will be assigned

55

- The resultant set of rules will probably not have the same coverage as the decision tree
- Its accuracy should be equivalent
- Rules are much easier to understand
- Rules can be tuned by hand by an expert

56

From Trees to Rules

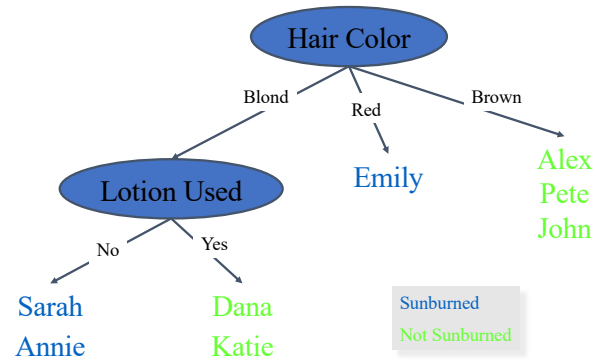
Once an identification tree is constructed, it is a simple matter to concert it into a set of equivalent rules

- Example from Artificial Intelligence, P.H. Winston 1992

Independent Attributes / Condition Attributes					Dependent Attributes / Decision Attributes
Name	Hair	Height	Weight	Lotion	Result
Sarah	blonde	average	light	no	sunburned (positive)
Dana	blonde	tall	average	yes	none (negative)
Alex	brown	short	average	yes	none
Annie	blonde	short	average	no	sunburned
Emily	red	average	heavy	no	sunburned
Pete	brown	tall	heavy	no	none
John	brown	average	heavy	no	none
Katie	blonde	short	light	yes	none

57

An ID3 tree consistent with the data



58

Corresponding rules

*If the person's hair is blonde
and the person uses lotion
then nothing happens*

*If the person's hair color is blonde
and the person uses no lotion
then the person turns red*

*If the person's hair color is red
then the person turns red*

*If the person's hair color is brown
then nothing happens*

59

Unnecessary Rule Antecedents should be eliminated

*If the person's hair is blonde
and the person uses lotion
then nothing happens*

Are both antecedents are really necessary?
Dropping the first antecedents produce a rule with the same results

*If the the person uses lotion
then nothing happens*

60

Unnecessary Rules should be Eliminated

*If the person uses lotion
then nothing happens*

*If the person's hair color is blonde
and the person uses no lotion
then the person turns red*

*If the person's hair color is red
then the person turns red*

*If the person's hair color is brown
then nothing happens*

61

- Note that two rules have a consequent that indicate that a person will turn red, and two that indicate that nothing happens
- One can replace either the two of them by a **default rule**

62

Default rule

*If the person uses lotion
then nothing happens*

*If the person's hair color is brown
then nothing happens*

***If no other rule applies
then the person turns red***

63

Continuous input variables

- Problem? Previous principles only applicable to discrete variables
 - how to handle **continuous variables**?
- Solutions:
 - **variable discretization**
 - e.g. income in the credit risk example
 - leave numeric values as-is and let the tree learning approach identify the **best binarization threshold (CART)**
 - **For example, the binarization of**
 - when selecting a continuous variable, examine possible split points for the real values
 - the **binary** split point that maximizes the discriminative power (information gain) is taken as a candidate

64

What is CART?

- Classification And Regression Trees
- Developed by Breiman, Friedman, Olshen, Stone in early 80's.
 - Introduced tree-based modeling into the statistical mainstream
 - Rigorous approach involving cross-validation to select the optimal tree
- One of many tree-based modeling techniques.
 - CART -- the classic
 - CHAID
 - C5.0
 - Software package variants (SAS, S-Plus, R...)
 - Note: the "rpart" package in "R" is freely available

65

Idea: Recursive Partitioning

- Take all of your data.
- Consider *all* possible values of *all* variables.
- Select the variable/value ($X=t_1$) that produces the greatest "separation" in the target.
 - ($X=t_1$) is called a "split".
- If $X < t_1$ then send the data to the "left"; otherwise, send data point to the "right".
- Now repeat same process on these two "nodes"
 - You get a "tree"
 - Note: CART only uses *binary* splits.

66

Split

- On the theoretical level, is very natural for a decision tree to handle categorical variables, scikit-learn don't do it and only accept continuous variables
- Scikit-learn **only** supports **binary splits** for **numerical values** ☹
 - Bigger trees

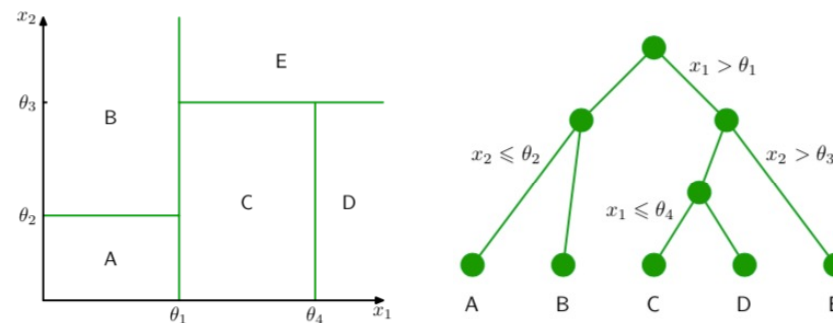
EXAMPLE:

- The best split in age is between 59 and 65
 - Only six possible splits since we can order the numbers..
- The best split in BMI is between 28 and 33
- Considering both splitting ranges, age has the highest information gain

age	BMI	hospitalization
33	17	Y
65	33	Y
68	35	Y
19	28	N
44	37	N
53	25	N
59	22	N

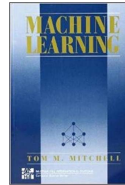
67

Binary Tree

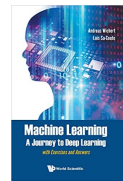


68

Literature



- Tom M. Mitchell, Machine Learning, McGraw-Hill; 1st edition (October 1, 1997)
 - Chapter 3



- Machine Learning - A Journey to Deep Learning, A. Wichert, Luis Sa-Couto, World Scientific, 2021
 - Chapter 2, Section 2.5
 - C) Decision Trees