



NUI Galway
OÉ Gaillimh

Topic 2: Information-based Learning

Part 1: Introduction



Learning objectives

After completing this topic successfully, you will be able to ...

1. Explain what supervised learning is
2. Distinguish it from unsupervised learning and reinforcement learning
3. Describe in detail an algorithm for decision tree induction
4. Apply decision tree induction to a data set
5. List related algorithms
6. Discuss high-level concepts such as choice of hypothesis language, overfitting, underfitting and noise

Reading: Russell & Norvig 3rd Ed, Chapter 18.18.4; Kelleher et al. Chapter 4



Overview of topic

This week:

1. Introduction, learning objectives and overview
2. Supervised learning principles
3. Decision trees
4. Entropy
5. Information gain

Next week:

6. The ID3 algorithm
7. Issues in decision tree learning
8. ID3 extensions and related algorithms
9. Supervised learning considerations
10. Review of topic



NUI Galway
OÉ Gaillimh

Topic 2: Information-based learning

Part 2: Supervised learning principles



Supervised learning: motivating examples

1. Estimate sale price of a house, given past data of house sizes, locations and their prices
2. Before unlocking a tablet, determine whether a known user or somebody else is looking at the webcam
3. Decide whether a chemical spectrum of a mixture has evidence of containing cocaine, based on other spectra with & without cocaine
4. Predict concentration of cocaine in mixture
5. Determine whether objects of interest are present in a scene – if so, what are they? (relevant for autonomous vehicles and robotics, among other domains)

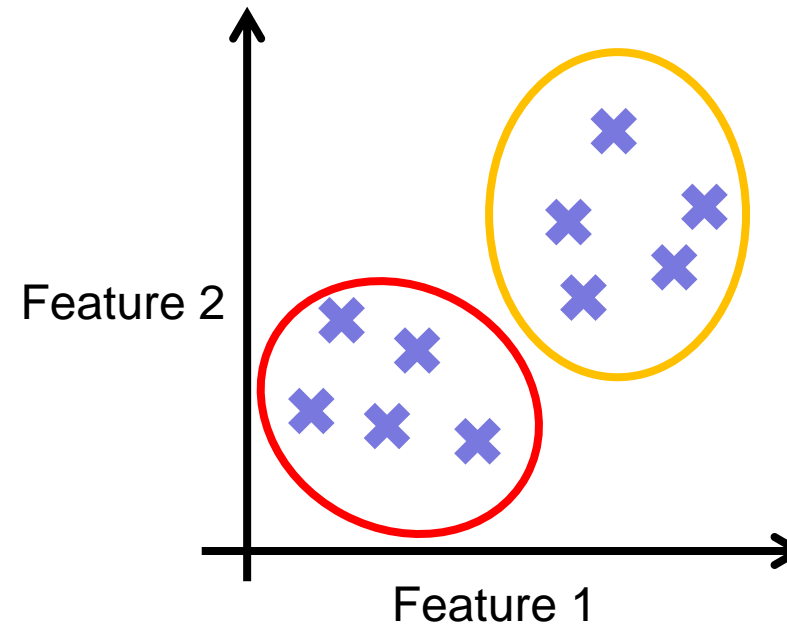
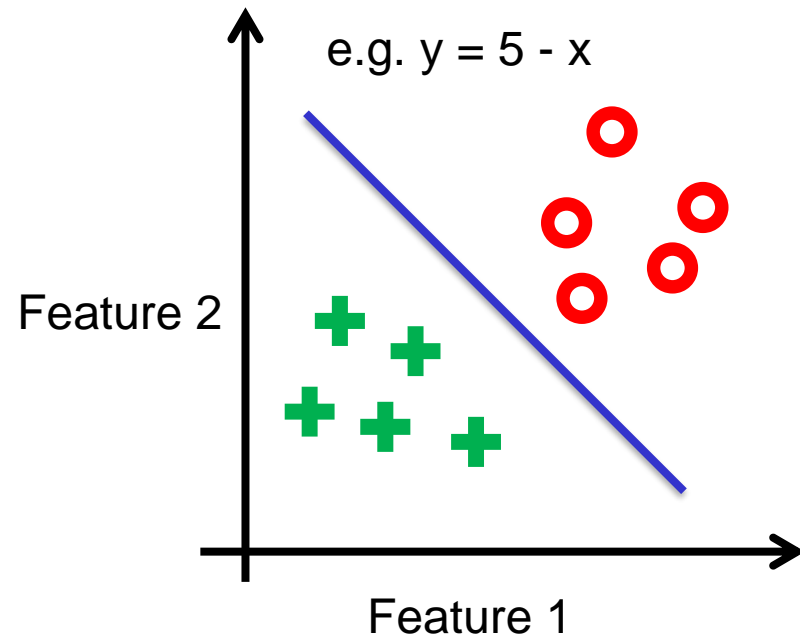
Key feature:

given "right answers" / "ground truth" as start point.

Which tasks are classification, and which are regression?



Supervised vs. unsupervised learning





Supervised learning: task definition

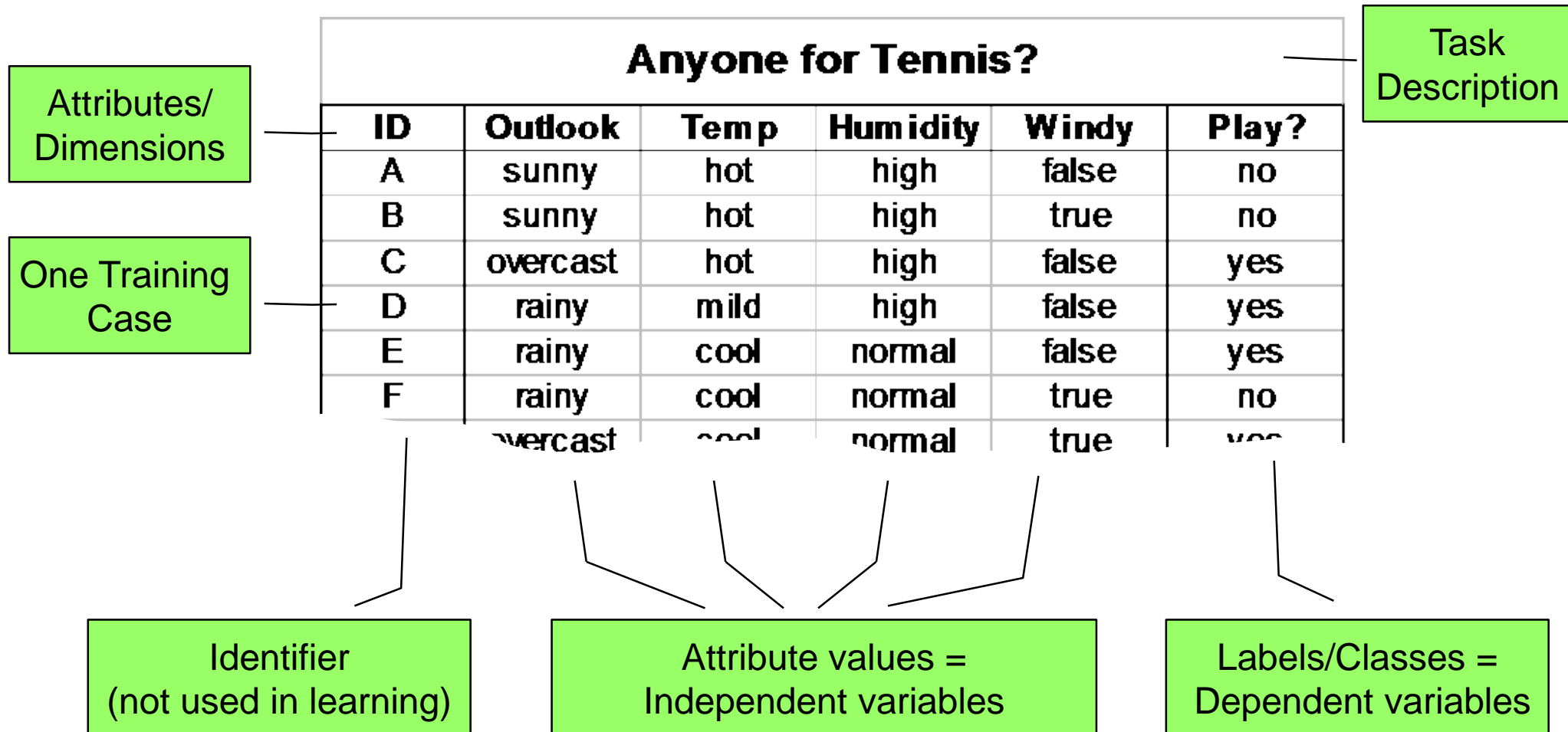
- Given examples, return function h (*hypothesis*) that approximates some 'true' function f that (hypothetically) generated the labels for the examples
 - Have set of examples, the **training data**:
each has a **label** and a set of **attributes** that have known **values**
 - Consider *labels* (classes) to be *outputs* of some function f ; the observed *attributes* are its *inputs*
 - Denote the attribute value inputs \mathbf{x} , labels are their corresponding outputs $f(\mathbf{x})$
 - An example is a pair $(\mathbf{x}, f(\mathbf{x}))$
 - Function f is *unknown*; want to discover an approximation of it, h
 - Can use h to **predict** labels of new data: **generalisation**

Also known as Pure Inductive Learning – **why?**



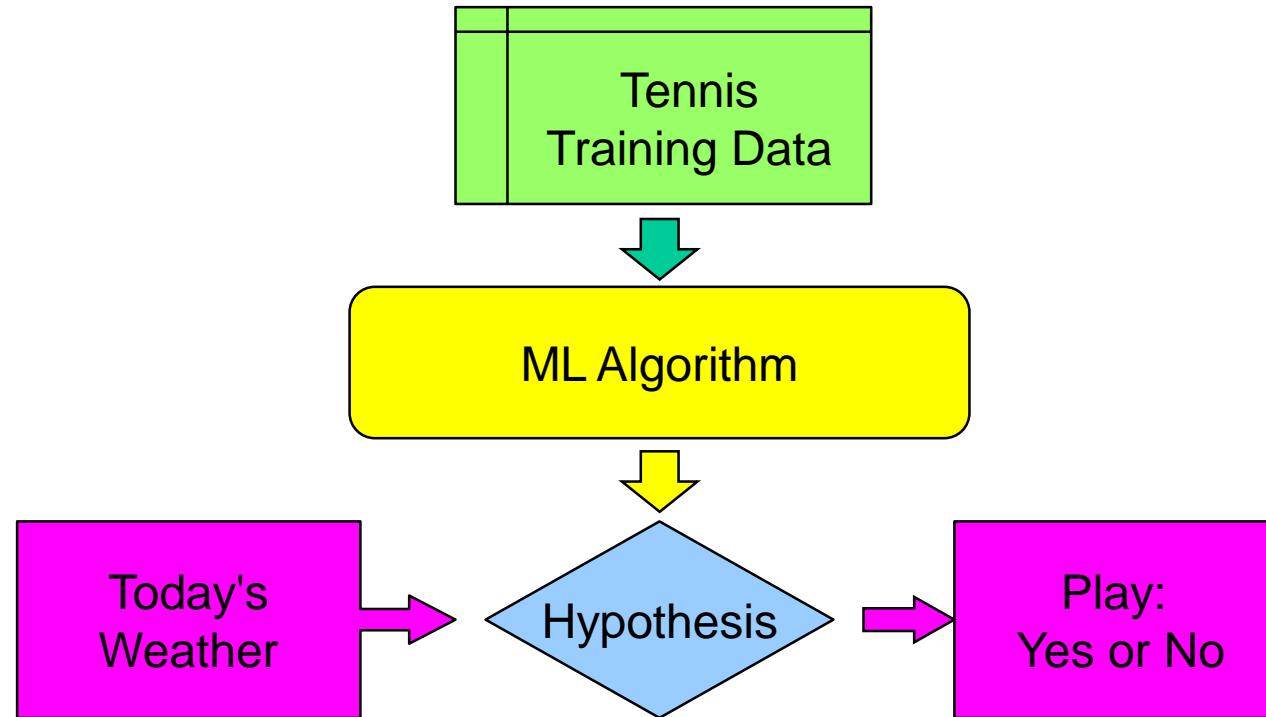


Training data example





Overview of the supervised learning process





NUI Galway
OÉ Gaillimh

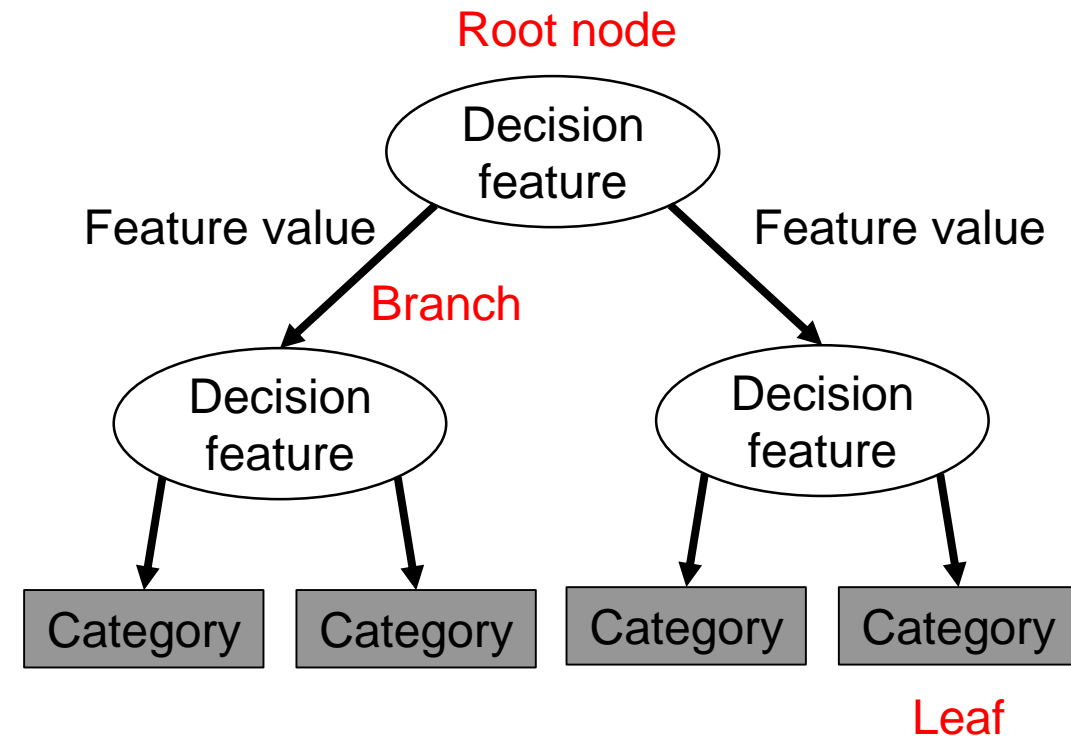
Topic 2: Information-based Learning

Part 3: Decision trees



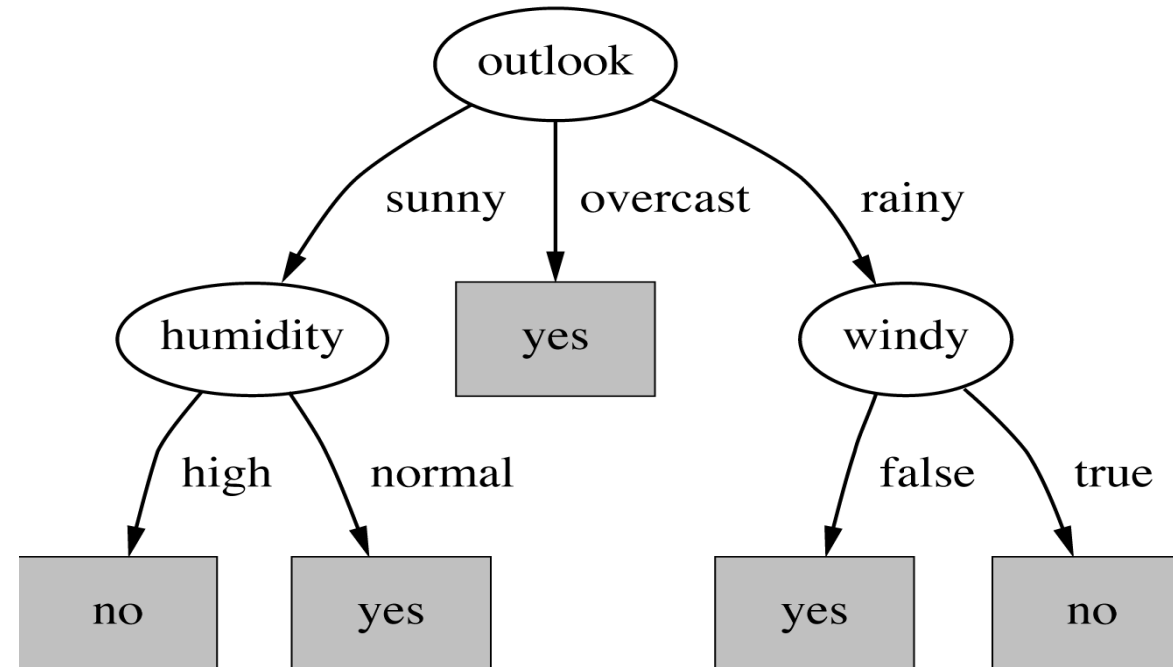
Decision trees

- Decision trees are a fundamental structure used in information-based machine learning
- Main idea: use a decision tree as a predictive model, to decide what category/label/class an item belongs to based on the values of its features
- So-called due to their tree-like structure:
 - A node (where two branches intersect) is a decision point. Nodes partition the data.
 - observations about an item (values of features) are represented using branches
 - The terminal nodes are called leaves; these specify the target label for an item





Decision tree for a sample dataset





Example dataset for induction (1)

- Weather dataset
 - Four attributes:
 - outlook**: sunny / overcast / rainy
 - temperature**: hot / mild / cool
 - humidity**: high / normal
 - windy**: true / false
 - Used to decide whether or not to *play tennis*
 - 14 examples in dataset
 - See **weather.xls** (spreadsheet) or **weathertext.csv** (comma separated values format)
- Objective:
 - Find hypothesis that *describes the cases* given and can be used to *make decisions* in other cases
 - Express the hypothesis as a decision tree.

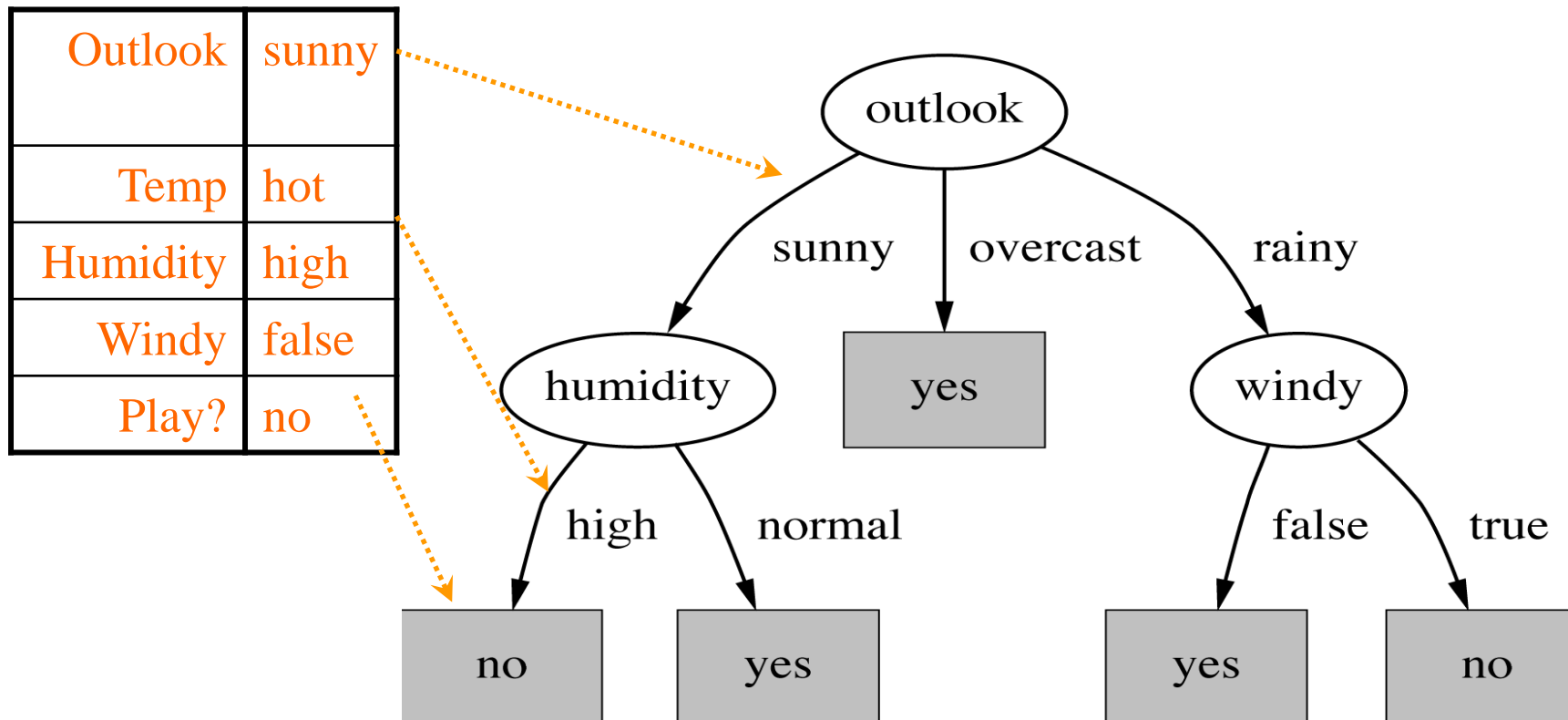


Example dataset for induction (2)

Anyone for Tennis?					
ID	Outlook	Temp	Humidity	Windy	Play?
A	sunny	hot	high	false	no
B	sunny	hot	high	true	no
C	overcast	hot	high	false	yes
D	rainy	mild	high	false	yes
E	rainy	cool	normal	false	yes
F	rainy	cool	normal	true	no
G	overcast	cool	normal	true	yes
H	sunny	mild	high	false	no
I	sunny	cool	normal	false	yes
J	rainy	mild	normal	false	yes
K	sunny	mild	normal	true	yes
L	overcast	mild	high	true	yes
M	overcast	hot	normal	false	yes
N	rainy	mild	high	true	no



Decision tree for this data (1)

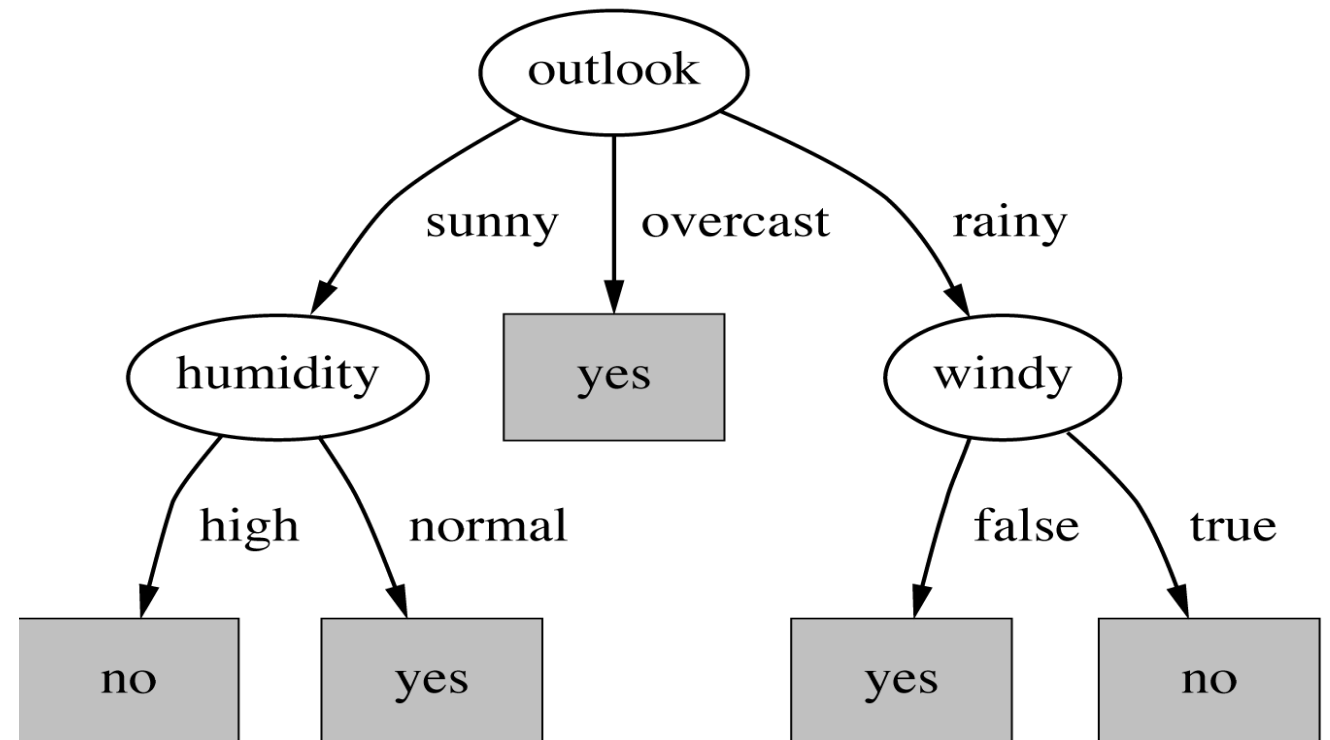




Decision tree for this data (2)

Anyone for Tennis?

ID	Outlook	Temp	Humidity	Windy	Play?
A	sunny	hot	high	false	no
B	sunny	hot	high	true	no
C	overcast	hot	high	false	yes
D	rainy	mild	high	false	yes
E	rainy	cool	normal	false	yes
F	rainy	cool	normal	true	no
G	overcast	cool	normal	true	yes
H	sunny	mild	high	false	no
I	sunny	cool	normal	false	yes
J	rainy	mild	normal	false	yes
K	sunny	mild	normal	true	yes
L	overcast	mild	high	true	yes
M	overcast	hot	normal	false	yes
N	rainy	mild	high	true	no





Inductive learning of a decision tree

Step 1

- For all attributes that have not yet been used in the tree, calculate their **entropy** and **information gain** values for the training samples

Step 2

- Select the attribute that has the highest information gain

Step 3

- Make a tree node containing that attribute

Repeat

- This node **partitions** the data:
apply the algorithm **recursively** to each partition



NUI Galway
OÉ Gaillimh

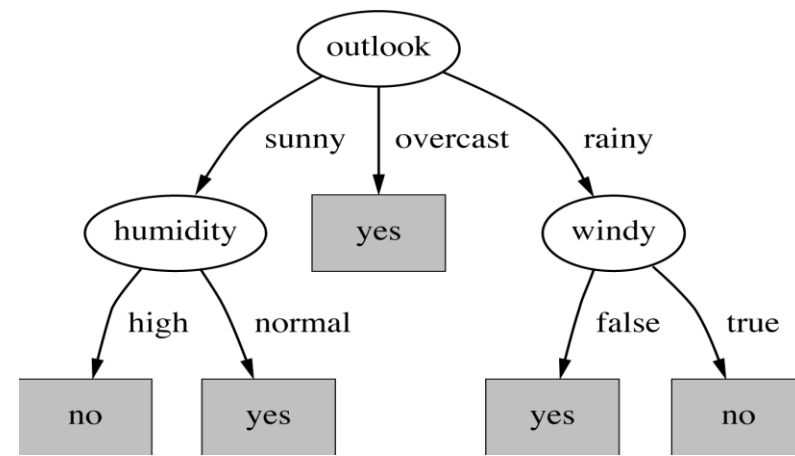
Topic 2: Information-based Learning

Part 4: Entropy



Motivation

- We already saw how some descriptive features can more effectively discriminate between (or predict) classes which are present in the dataset
- Decision trees partition the data at each node, so it makes sense to use features which have higher discriminatory power “higher up” in a decision tree.
- Therefore we need to develop a formal measure of the discriminatory power of a given attribute
- **Information gain – this can be calculated using entropy**

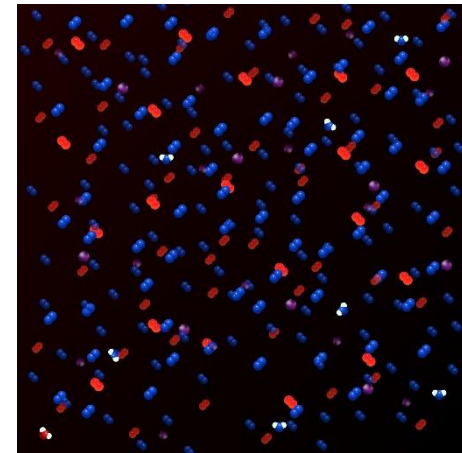


Anyone for Tennis?					
ID	Outlook	Temp	Humidity	Windy	Play?
A	sunny	hot	high	false	no
B	sunny	hot	high	true	no
C	overcast	hot	high	false	yes
D	rainy	mild	high	false	yes
E	rainy	cool	normal	false	yes
F	rainy	cool	normal	true	no
G	overcast	cool	normal	true	yes
H	sunny	mild	high	false	no
I	sunny	cool	normal	false	yes
J	rainy	mild	normal	false	yes
K	sunny	mild	normal	true	yes
L	overcast	mild	high	true	yes
M	overcast	hot	normal	false	yes
N	rainy	mild	high	true	no



Entropy

- Claude Shannon (often referred to as “the father of information theory”) proposed a measure to of the impurity of the elements in a set, referred to as entropy
- Entropy may be used to measure of the uncertainty of a random variable
- The term entropy generally refers to disorder or uncertainty, so the use of this term in the context of information theory is analogous to the other well-known use of the term in statistical thermodynamics
- Acquisition of information (information gain) corresponds to a reduction in entropy
- “Information is the resolution of uncertainty” (Shannon)
- 1948 article “A Mathematical Theory of Communication”





Calculating entropy

- The entropy of a dataset S with n different classes may be calculated as:

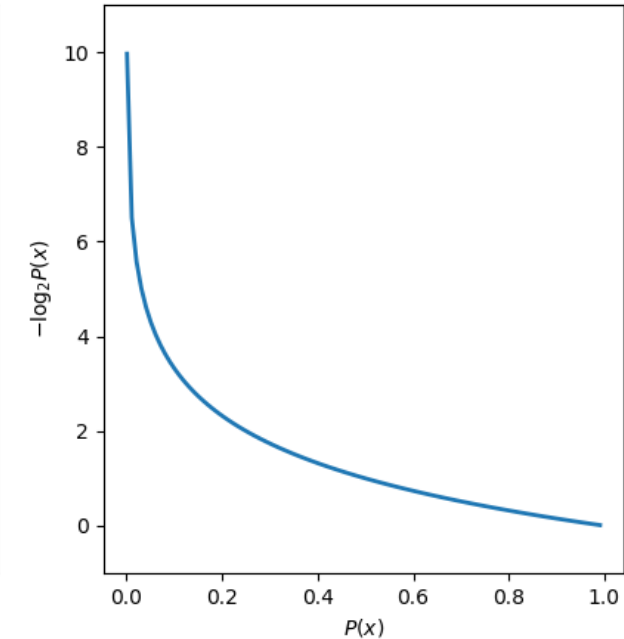
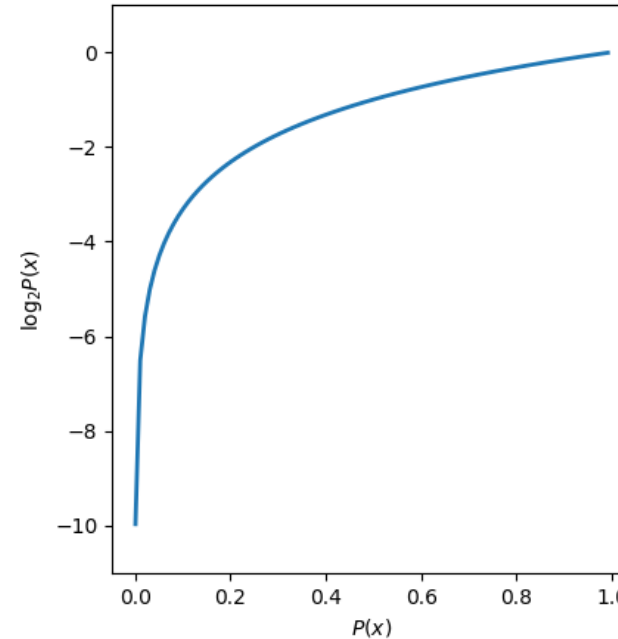
$$\text{Ent}(S) = \sum_{i=1}^n -p_i \log_2 p_i$$

- Here p_i is the proportion of class i in the dataset.
- This is an example of a probability mass function
- Entropy is typically measured in bits (note \log_2 in the equation above)
- The lowest possible entropy output from this function is 0 ($\log_2 1 = 0$)
- The highest possible entropy is $\log_2 n$ ($=1$ when there are only 2 classes)



Why use the binary logarithm?

- A useful measure of uncertainty should:
 - Assign high uncertainty values to outcomes with a low probability
 - Assign low uncertainty values to outcomes with a high probability
- Consider the plot to the right
 - \log_2 returns large negative values when P is close to 0
 - \log_2 returns small negative values when P is close to 1
- Using $-\log_2$ is more convenient, as this will give positive entropy values, with 0 as the lowest entropy





Entropy worked example 1

$$\text{Ent}(S) = \sum_{i=1}^n -p_i \log_2 p_i$$

$$\text{Ent}(S) = \text{Ent}([9+,5-])$$

$$\text{Ent}(S) = -9/14 \log_2(9/14) - 5/14 \log_2(5/14)$$

$$\text{Ent}(S) = 0.9403$$

If calculating this in a spreadsheet application such as Excel, make sure that you are using \log_2 (e.g. $\text{LOG}(9/14, 2)$)

Anyone for Tennis?					
ID	Outlook	Temp	Humidity	Windy	Play?
A	sunny	hot	high	false	no
B	sunny	hot	high	true	no
C	overcast	hot	high	false	yes
D	rainy	mild	high	false	yes
E	rainy	cool	normal	false	yes
F	rainy	cool	normal	true	no
G	overcast	cool	normal	true	yes
H	sunny	mild	high	false	no
I	sunny	cool	normal	false	yes
J	rainy	mild	normal	false	yes
K	sunny	mild	normal	true	yes
L	overcast	mild	high	true	yes
M	overcast	hot	normal	false	yes
N	rainy	mild	high	true	no



Entropy worked example 2

Anyone for Tennis?

ID	Outlook	Temp	Humidity	Windy	Play?
A	sunny	hot	high	false	no
B	sunny	hot	high	true	no
C	overcast	hot	high	false	yes
D	rainy	mild	high	false	yes
E	rainy	cool	normal	false	yes
F	rainy	cool	normal	true	no
G	overcast	cool	normal	true	yes
H	sunny	mild	high	false	no
I	sunny	cool	normal	false	yes
J	rainy	mild	normal	false	yes
K	sunny	mild	normal	true	yes
L	overcast	mild	high	true	yes
M	overcast	hot	normal	false	yes
N	rainy	mild	high	true	no

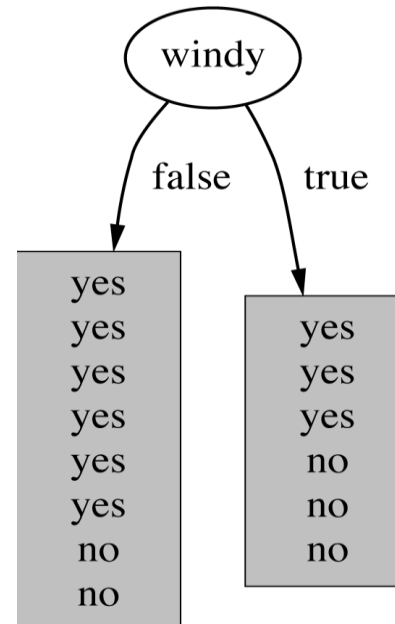


Entropy worked example 2

$$\text{Ent}(S) = \sum_{i=1}^n -p_i \log_2 p_i$$

$$\begin{aligned}\text{Ent}(S_{\text{windy=false}}) &= \text{Ent}([6+, 2-]) \\ &= -6/8 \log_2(6/8) - 2/8 \log_2(2/8) \\ &= 0.3112 + 0.5 = 0.8112\end{aligned}$$

$$\begin{aligned}\text{Ent}(S_{\text{windy=true}}) &= \text{Ent}([3+, 3-]) \\ &= -3/6 \log_2(3/6) - 3/6 \log_2(3/6) \\ &= 0.5 + 0.5 = \mathbf{1.0}\end{aligned}$$



Anyone for Tennis?					
ID	Outlook	Temp	Humidity	Windy	Play?
A	sunny	hot	high	false	no
B	sunny	hot	high	true	no
C	overcast	hot	high	false	yes
D	rainy	mild	high	false	yes
E	rainy	cool	normal	false	yes
F	rainy	cool	normal	true	no
G	overcast	cool	normal	true	yes
H	sunny	mild	high	false	no
I	sunny	cool	normal	false	yes
J	rainy	mild	normal	false	yes
K	sunny	mild	normal	true	yes
L	overcast	mild	high	true	yes
M	overcast	hot	normal	false	yes
N	rainy	mild	high	true	no



NUI Galway
OÉ Gaillimh

Topic 2: Information-based Learning

Part 5: Information gain



Information gain

- The **information gain** of an attribute is the reduction in entropy from partitioning the data according to that attribute

$$\text{Gain}(S, A) = \text{Ent}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Ent}(S_v)$$

- Here S is the entire set of data being considered, and S_v refers to each partition of the data according to each possible value v for the attribute
- $|S|$ and $|S_v|$ refer to the cardinality or size of the overall dataset, and the cardinality or size of a partition respectively
- When selecting an attribute for a node in a decision tree, use whichever attribute A gives the greatest information gain



Information gain worked example

$$\text{Gain}(S, A) = \text{Ent}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Ent}(S_v)$$

$$|S| = 14$$

$$|S_{\text{windy}=\text{true}}| = 6$$

$$|S_{\text{windy}=\text{false}}| = 8$$

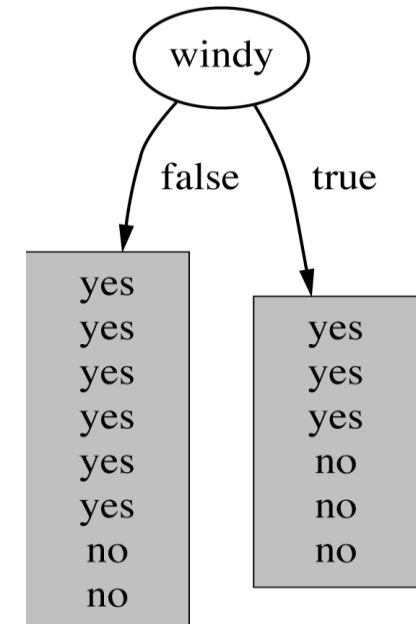
$$\text{Gain}(S, \text{Windy})$$

$$= \text{Ent}(S) - |S_{\text{windy}=\text{true}}|/|S| \text{Ent}(S_{\text{windy}=\text{true}}) - |S_{\text{windy}=\text{false}}|/|S| \text{Ent}(S_{\text{windy}=\text{false}})$$

$$= \text{Ent}(S) - (6/14) \text{Ent}([3+, 3-]) - (8/14) \text{Ent}([6+, 2-])$$

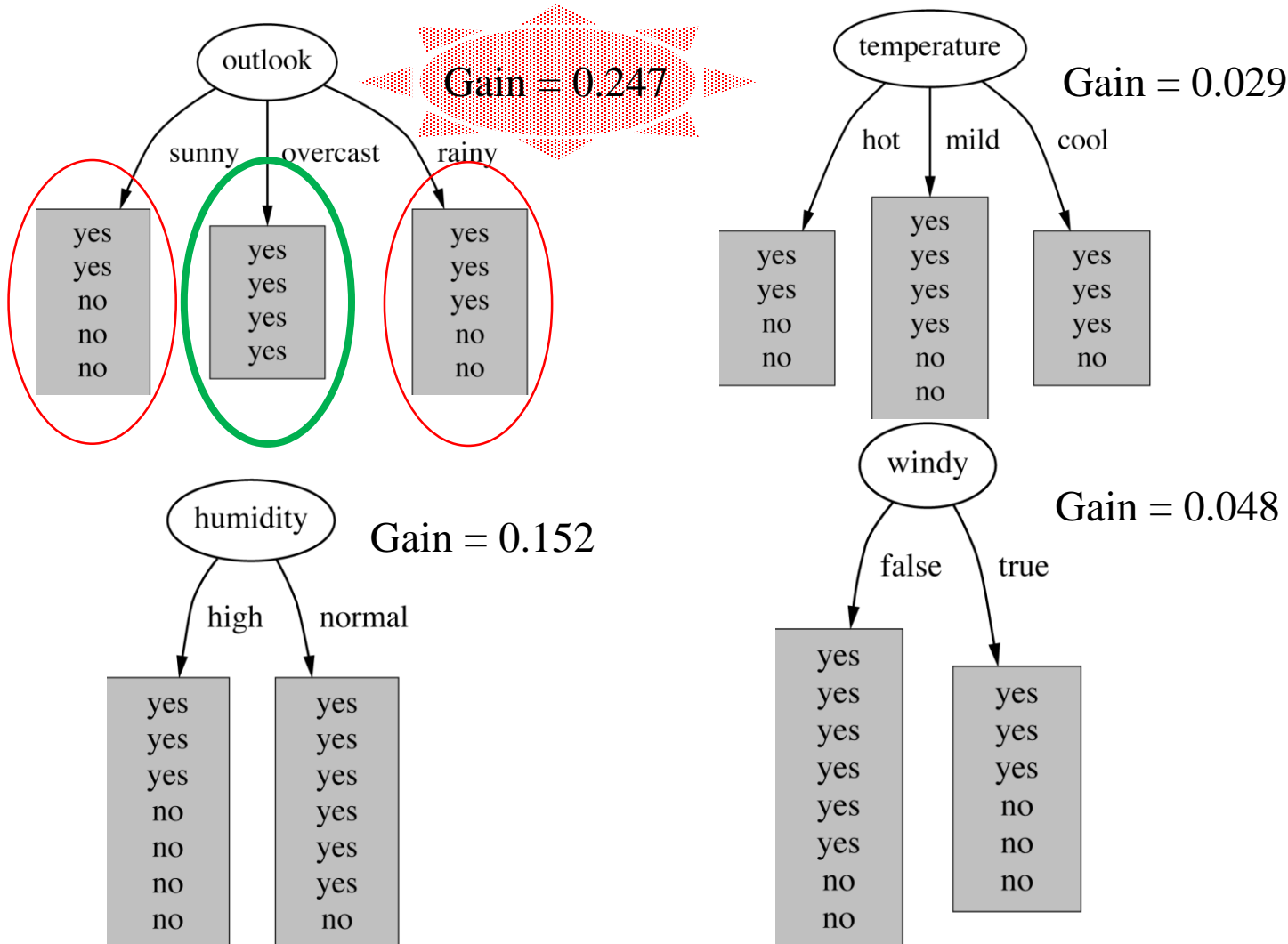
$$= 0.940 - (6/14) 1.00 - (8/14) 0.811$$

$$\text{Gain}(S, \text{Windy}) = \mathbf{0.048}$$





Best partitioning = highest information gain



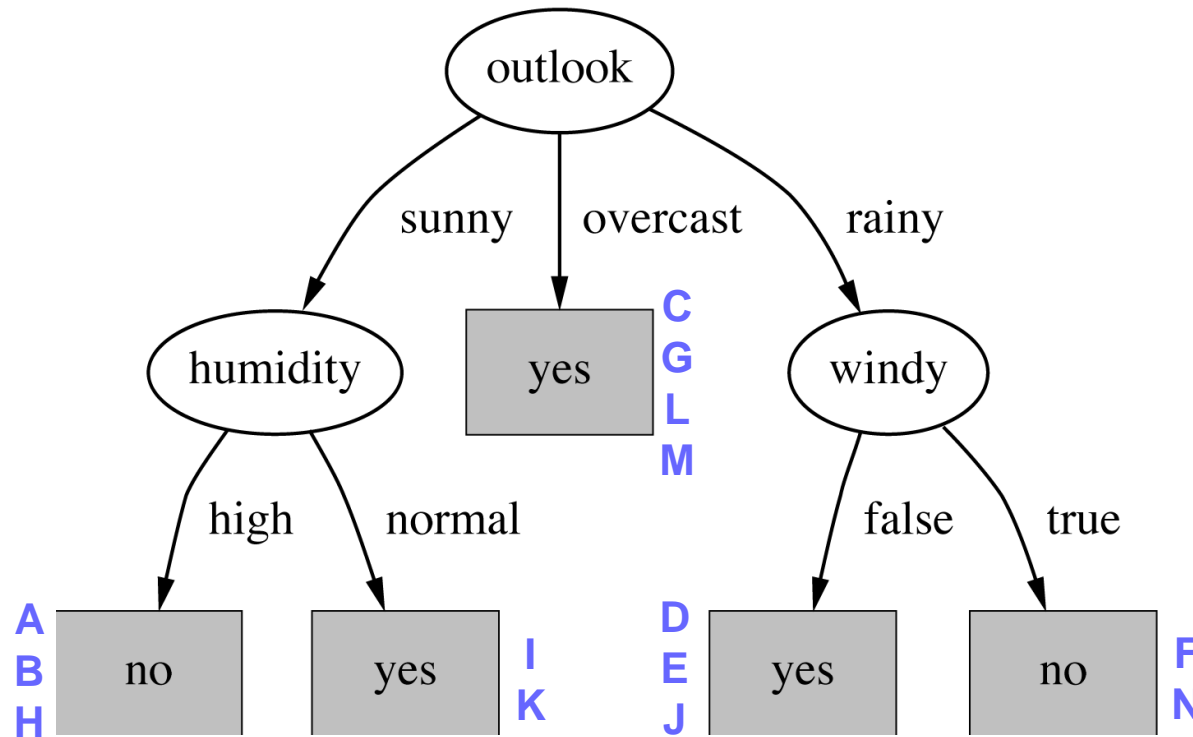
Anyone for Tennis?					
ID	Outlook	Temp	Humidity	Windy	Play?
A	sunny	hot	high	false	no
B	sunny	hot	high	true	no
C	overcast	hot	high	false	yes
D	rainy	mild	high	false	yes
E	rainy	cool	normal	false	yes
F	rainy	cool	normal	true	no
G	overcast	cool	normal	true	yes
H	sunny	mild	high	false	no
I	sunny	cool	normal	false	yes
J	rainy	mild	normal	false	yes
K	sunny	mild	normal	true	yes
L	overcast	mild	high	true	yes
M	overcast	hot	normal	false	yes
N	rainy	mild	high	true	no

Having found the best split for the root node, repeat the whole procedure with each subset of examples ...

S will now refer to the subset in the partition being considered, instead of the entire dataset



Example: complete decision tree



Anyone for Tennis?					
ID	Outlook	Temp	Humidity	Windy	Play?
A	sunny	hot	high	false	no
B	sunny	hot	high	true	no
C	overcast	hot	high	false	yes
D	rainy	mild	high	false	yes
E	rainy	cool	normal	false	yes
F	rainy	cool	normal	true	no
G	overcast	cool	normal	true	yes
H	sunny	mild	high	false	no
I	sunny	cool	normal	false	yes
J	rainy	mild	normal	false	yes
K	sunny	mild	normal	true	yes
L	overcast	mild	high	true	yes
M	overcast	hot	normal	false	yes
N	rainy	mild	high	true	no

What about Temp = {Hot, Mild, Cool}?



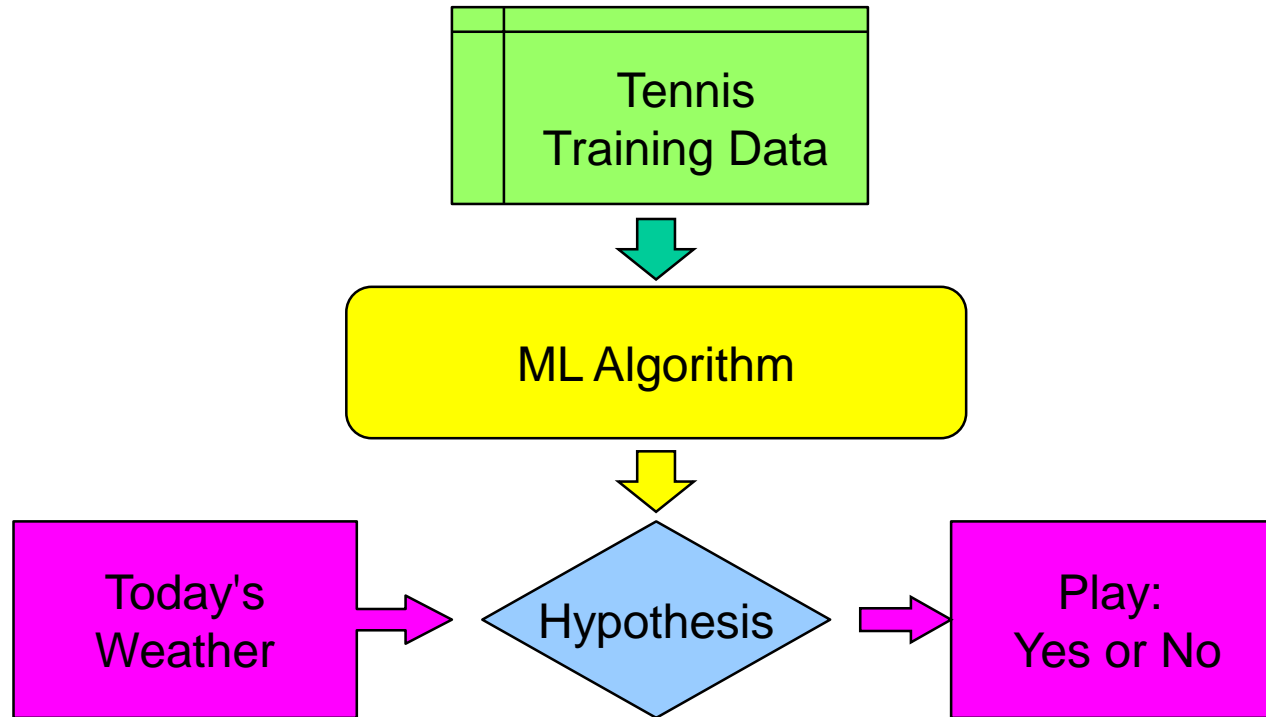
NUI Galway
OÉ Gaillimh

Topic 2: Information-based Learning

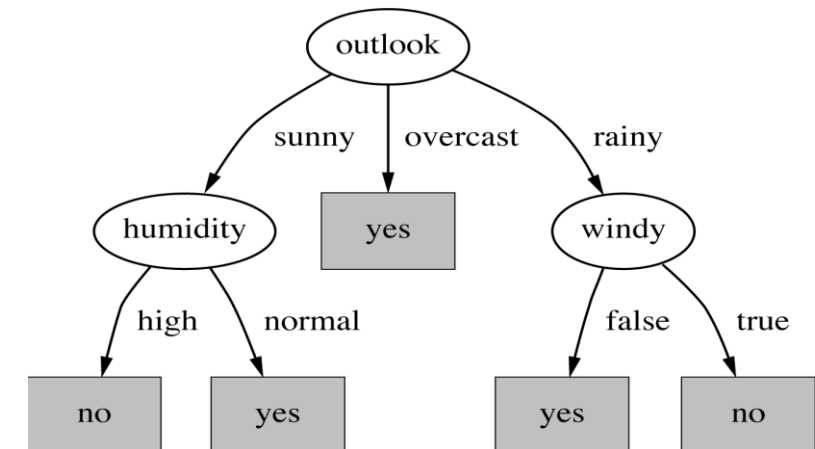
Part 6: The ID3 algorithm



Review: the supervised learning process



Anyone for Tennis?					
ID	Outlook	Temp	Humidity	Windy	Play?
A	sunny	hot	high	false	no
B	sunny	hot	high	true	no
C	overcast	hot	high	false	yes
D	rainy	mild	high	false	yes
E	rainy	cool	normal	false	yes
F	rainy	cool	normal	true	no
G	overcast	cool	normal	true	yes
H	sunny	mild	high	false	no
I	sunny	cool	normal	false	yes
J	rainy	mild	normal	false	yes
K	sunny	mild	normal	true	yes
L	overcast	mild	high	true	yes
M	overcast	hot	normal	false	yes
N	rainy	mild	high	true	no





Review: inductive learning of a decision tree

Step 1

- For all attributes that have not yet been used in the tree, calculate their **entropy** and **information gain** values for the training samples

Step 2

- Select the attribute that has the highest information gain

Step 3

- Make a tree node containing that attribute

Repeat

- This node **partitions** the data:
apply the algorithm **recursively** to each partition



The ID3 algorithm

1. ID3(Examples, Attributes, Target):
2. Input: **Examples**: set of classified examples
3. **Attributes**: set of attributes in the examples
4. **Target**: classification to be predicted
5. if **Examples** is empty then return a **Default** class
6. else if all **Examples** have same class then return this class
7. else if all **Attributes** are tested then return majority class
8. else:
9. let **Best** = attribute that **best separates** **Examples** relative to **Target**
10. let **Tree** = new decision tree with **Best** as root node
11. foreach value v_i of **Best**:
12. let **Examples_i** = subset of **Examples** that have **Best**= v_i
13. let **Subtree** = ID3(Examples_i, Attributes-Best, Target)
14. add branch from **Tree** to **Subtree** with label v_i
15. return **Tree**

Ross Quinlan, 1986



BASE
CASES

RECURSIVE
CALL

Based on
Algorithm
**Decision-Tree-
Learning** in
Russell &
Norvig textbook



NUI Galway
OÉ Gaillimh

Topic 2: Information-based Learning

Part 7: Issues in decision tree learning



Decision tree characteristics

- Popular because:
 - Relatively **easy** algorithm
 - **Fast**: greedy search without backtracking
 - **Comprehensible** output:
important in decision-making (medical, financial, ...)
 - **Practical**: discrete/numeric, irrelevant attributes, noisy data, ...
- Expressiveness: what functions can a DT represent?
 - Technically, any Boolean function (propositional logic)
 - Some functions, however, require exponentially large tree (e.g. parity function)
 - Cannot consider relationships between two attributes



Dealing with noisy or missing data

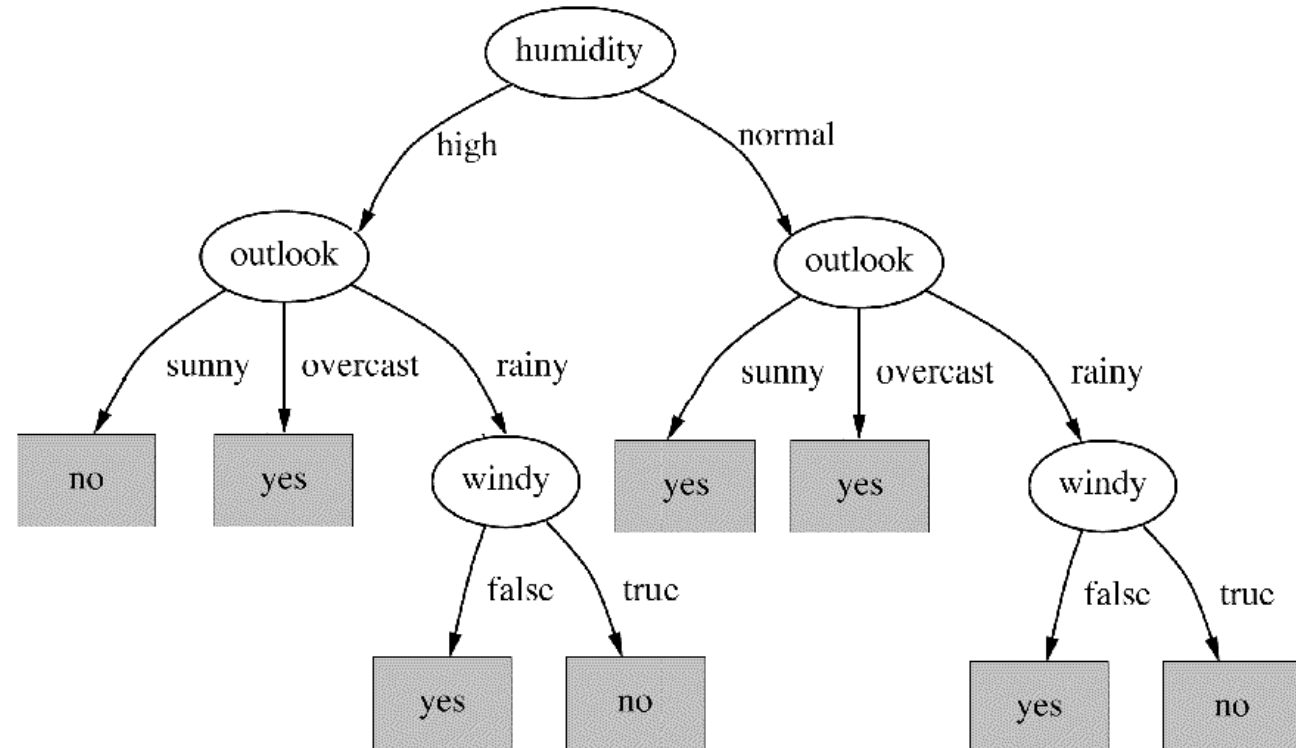
- What about inconsistent ("noisy") data?
 - Use majority class (line 7 of ID3 alg.)
 - 7. else if all **Attributes** are tested then return majority class
 - or interpret as probabilities
 - or return “average” target feature value
- What about missing data?
 - Given a complete decision tree, how should one classify an example that is missing one of the test attributes?
 - How should one modify the information gain formula when some training examples have unknown values for an attribute?
 - Could assign the most common value among the training examples that reach that node
 - Or could assume the attribute has all possible values, weighting each value according to its frequency among the training examples that reach that node



Instability of decision trees

- Hypothesis found is sensitive to training set used
 - consequence of greedy search
- Replace one example:
 - new one **consistent with original tree**
- Some algorithmic modifications to reduce the instability of decision tree learning were proposed by Li and Belford in their 2002 paper “Instability of decision tree classification algorithms”.
- Li and Belford’s main idea is to alter the attribute selection procedure, so that the tree learning algorithm is less sensitive to some % of the training dataset being replaced.

ID	Outlook	Temp	Humidity	Windy	Play?
C	overcast	hot	high	false	yes
O	sunny	hot	normal	true	yes





Pruning

- Overfitting occurs in a predictive model when the hypothesis learned makes predictions which are based on spurious patterns in the training data set.
Consequence: poor generalisation to new examples.
- Overfitting may happen for a number of reasons, including sampling variance or noise present in the dataset
- **Tree pruning** may be used to combat overfitting in decision trees
- Tree pruning can lead to induced trees which are inconsistent with the training set
- Generally, there are two different approaches to pruning:
 - Pre-pruning (e.g. \leq target # of examples in a partition, limiting tree depth, creating a new node only when information gain is above a threshold, statistical tests such as χ^2)
 - Post-pruning (e.g. target # of examples, compare error rate for model on a validation dataset with and without a given subtree; only keep a subtree if it improves the error rate, statistical tests such as χ^2 , reduced error pruning (Quinlan, 1987))



NUI Galway
OÉ Gaillimh

Topic 2: Information-based Learning

Part 7: ID3 extensions and related algorithms



Continuous-valued attributes

- What about continuous-valued attributes?
 - Pick threshold value T for attribute A , and test whether $A > T$
 - Information Gain can be used to decide which T is best
 - Could select T at a **midpoint** where classification changes

Temp	40	48	54	60	72	80	85	90
Play?	No	No		Yes	Yes	Yes		No



Selecting the best attribute: alternative metrics (1)

- Earlier, we introduced the concept of information gain, which we can use as a metric for the discriminatory power of an attribute
- Information gain does have some drawbacks; it tends to favour attributes that can take on a large number of different values
- One alternative is to use the **information gain ratio**

$$\text{GainRatio}(S, A) = \frac{\text{Gain}(S, A)}{\sum_{v \in \text{Values}(A)} -p_v \log_2 p_v}$$

- The divisor of this fraction measures the amount of information used to compute the gain value, and is the entropy of S with respect to A



Selecting the best attribute: alternative metrics (2)

- Another alternative is to use the **Gini index** instead of entropy as a measure of the impurity of a set

$$\text{Gini}(S) = 1 - \sum_{i=1}^n p_i^2$$

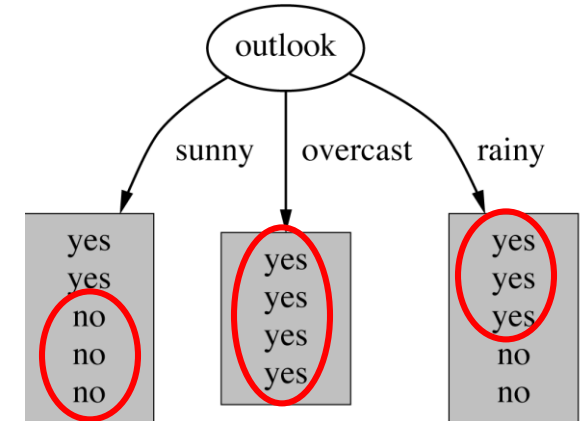
- Then the gain for a feature may be calculated based on the reduction in the Gini index (rather than a reduction in entropy):

$$\text{GiniGain}(S, A) = \text{Gini}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Gini}(S_v)$$



Related algorithms [1]

- 1R
 - Decision tree with just one rule
 - Introduced in a paper by Robert C. Holte (1993).
“Very Simple Classification Rules Perform Well on Most Commonly Used Datasets”, Computer Science Department, University of Ottawa



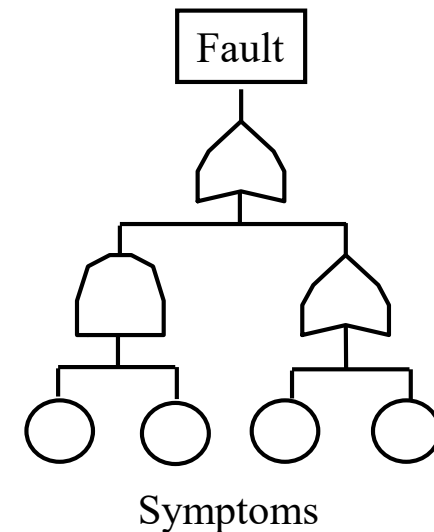
- Decision Stump
 - 1 rule with 1 test
 - Outlook = Overcast => YES
 - Outlook != Overcast => YES

These are deliberately simple variants that are used within other algorithms (meta-learning; ensembles). Often referred to as “weak learners”.



Related algorithms [2]

- Decision Lists
 - A set of rules (predicate logic), describing the hypothesis, that are followed in the given order
- C4.5 Rules
 - Alternative representation of C4.5 decision trees
- PART: Rules constructed with *partial* DTs
 - Can be more readable than standard DTs
- IFT: Induction of Fault Trees





Decision tree software

- C4.5
 - Original implementation (C language, command line), available on Ross Quinlan's website (<https://www.rulequest.com/Personal/>)
 - Deals with missing values in the data, high-branching attributes (e.g. ID in the weather data), pruning to avoid overfitting, converting a decision tree to a list of rules
- C5.0
 - Commercial version from RuleQuest Research, with improvements over C4.5
- WEKA software
 - Accompanies book by Witten & Frank, Data Mining: Practical Machine Learning Tools and Techniques
 - Java implementations of many ML algorithms, including C4.5 (mysteriously called J48)
 - Easy-to-use front end and utilities
- Many other implementations in Python and R...



NUI Galway
OÉ Gaillimh

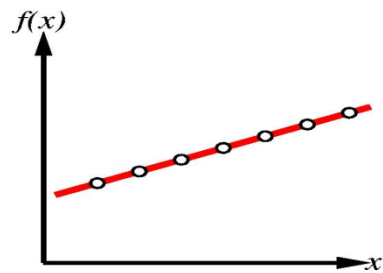
Topic 2: Information-based Learning

Part 8: Supervised learning considerations

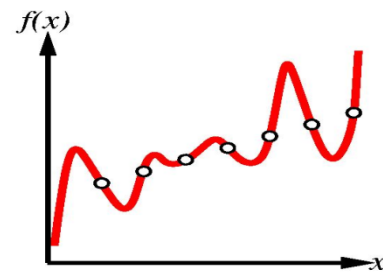


Supervised Learning Considerations [1]

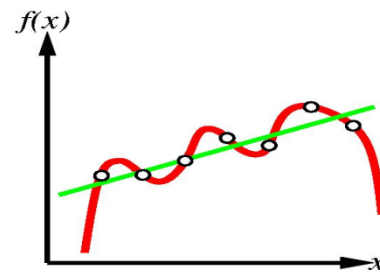
- Various hypotheses can be consistent with observations but inconsistent with each other:
Which one should we choose?



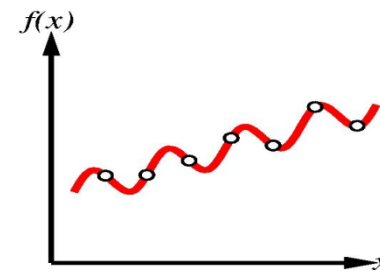
(1) Data with exact linear hypothesis



(2) Same data:
Exact 7th-Order
polynomial hyp.



(3) Different data:
Exact 5th-order poly
and approx. linear hyp.

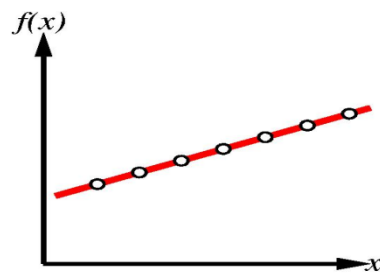


(4) Same data:
Exact sinusoidal
hypothesis

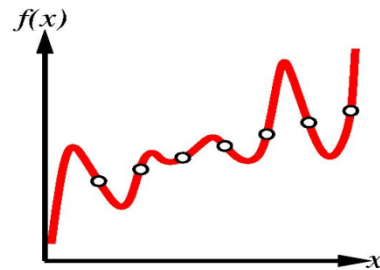


Supervised Learning Considerations [2]

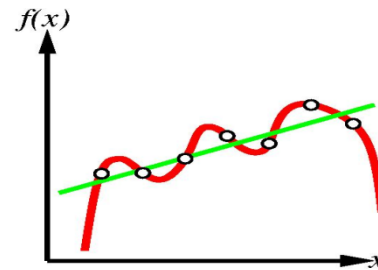
- Various hypotheses can be consistent with observations but inconsistent with each other:
Which one should we choose?



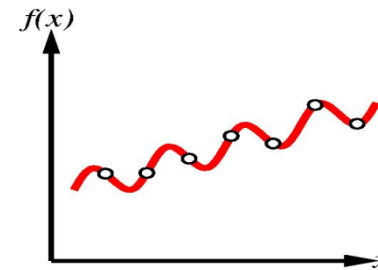
(1) Data with exact linear hypothesis



(2) Same data:
Exact 7th-Order
polynomial hyp.



(3) Different data:
Exact 5th-order poly
and approx. linear hyp.



(4) Same data:
Exact sinusoidal
hypothesis

- One solution: Ockham's Razor:
 - Prefer *simplest* hypothesis consistent with data
 - Definitions of simplicity (& consistency) may subject to debate
 - Depends strongly on how hypotheses are expressed



Supervised Learning Considerations [3]

- Hypothesis language is too limited?
 - Might be **unable** to find hypothesis that exactly matches 'true' function
 - If true function is more complex than what hypothesis can express, it will **underfit** the data
 - Saw this in previous slide, 3rd and 4th figures
- Hypothesis language cannot exactly match true function?
 - there will be a trade-off between **complexity** of hypothesis and how well it **fits the data**



Supervised Learning Considerations [4]

- Hypothesis language is very **expressive**?
 - Its search space is very large and the **computational complexity** of finding a good hypothesis will be high
 - Also need a large amount of data to avoid **overfitting**
- What can decision trees express?
 - Will learn about other algorithms that express hypotheses differently
 - In general, would like to use an algorithm for a problem that can express the true underlying function



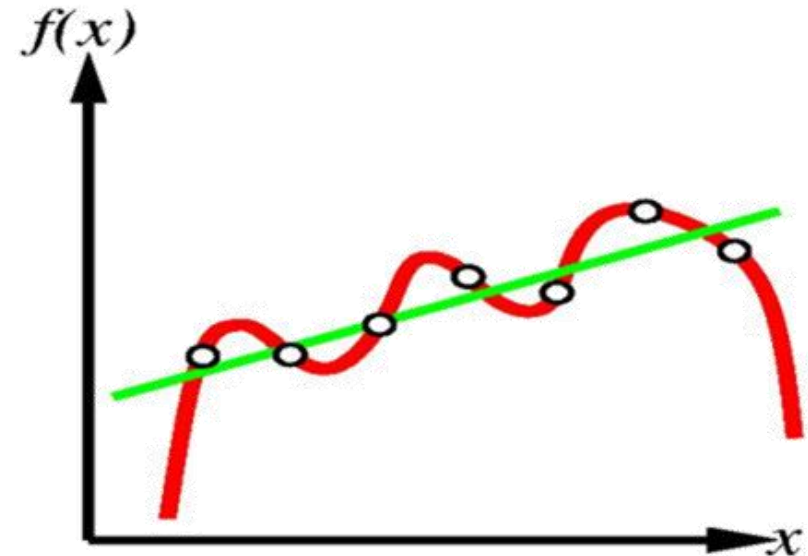
Supervised Learning Considerations [5]

- But don't forget:
we never know the true underlying function
- E.g. To avoid problem with poorly fitting data from a previous slide
 - Could change algorithm so that, as well as searching for coefficients of polynomials, it tries combinations of trig. functions (sin, cos, tan)
 - Learning problem will become enormously more complex, **but will it solve our problems?**
 - Probably not: you could easily think up some different kind of mathematical function, to generate a new dataset that the algorithm still cannot represent perfectly.
- For this reason, often use relatively simple hypothesis languages, in the absence of special knowledge about domain
 - more complex languages don't come with any real guarantees
 - more simple languages correspond to easier searching.



Noise, Overfitting and Underfitting [1]

- NOISE:
imprecise or incorrect attribute values or labels
 - Can't always quantify it, but should know from situation if it is present
 - E.g. labels may require subjective judgement or values may come from imprecise measurements





Noise, Overfitting and Underfitting [2]

- If the data might have noise, harder to decide which hypothesis is best:
 - Linear hypothesis could not fit to it, but polynomial could
 - But which would really be the better choice?
- Complex classification methods prone to **overfitting**;
simple ones prone to **underfitting**
 - If you increase complexity of hypothesis, you increase ability to fit to the data, but might also increase risk of overfitting

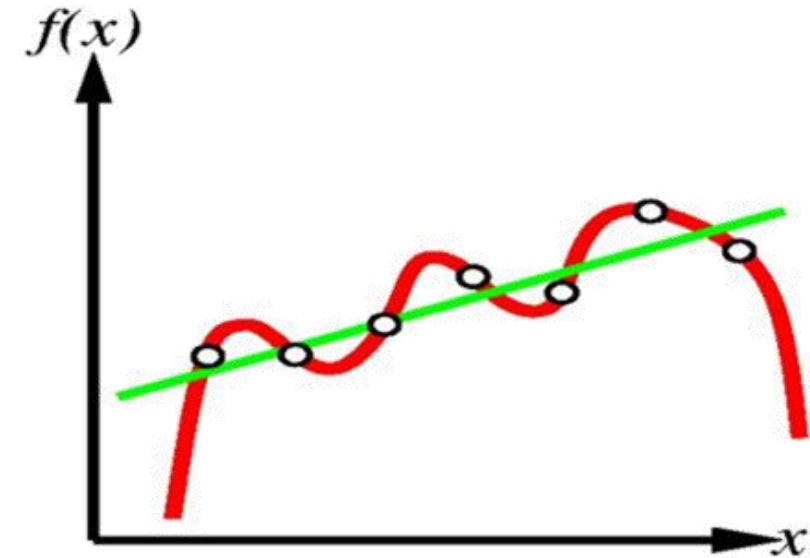
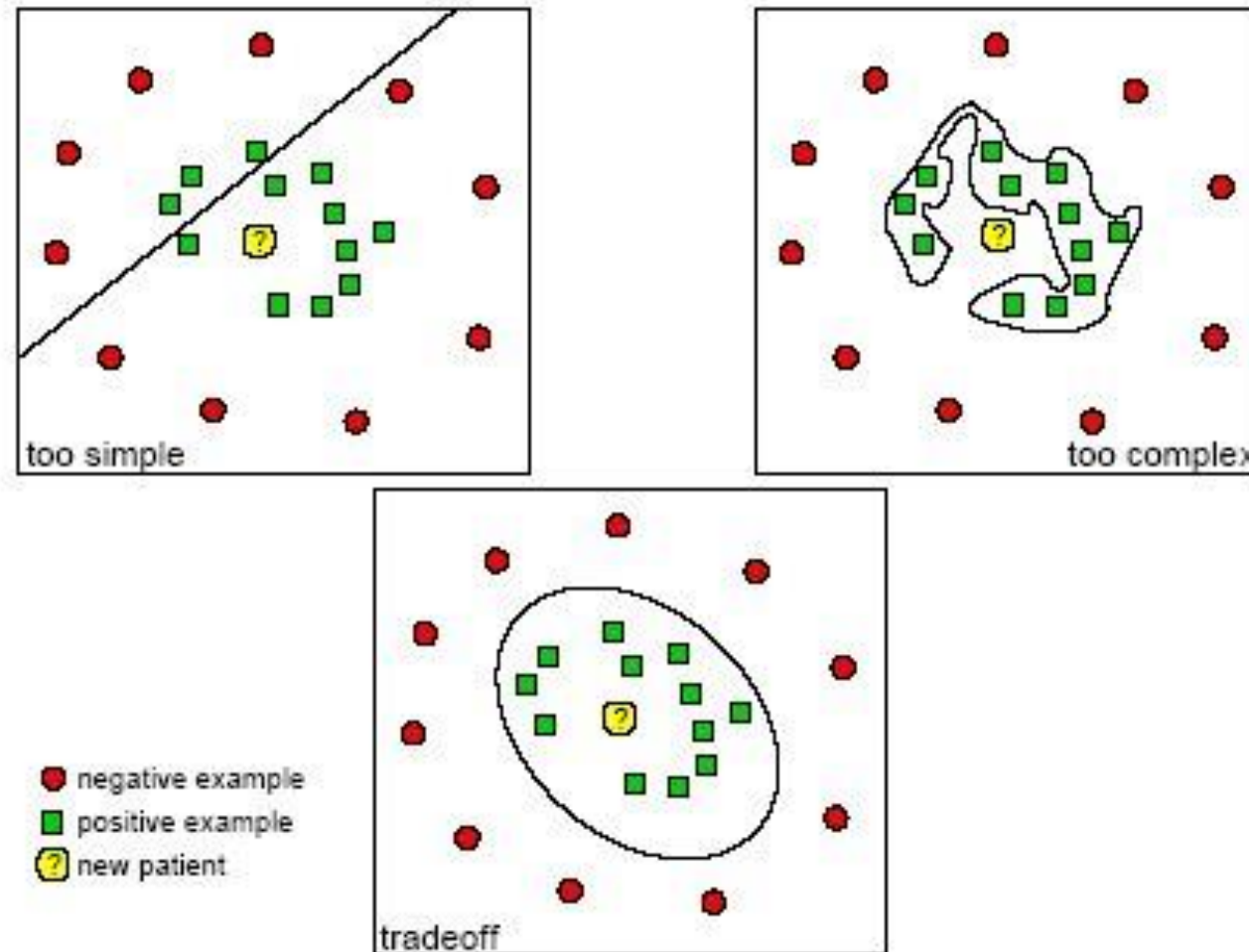




Illustration of Underfitting & Overfitting





Detecting Underfitting & Overfitting

- Previous slides have illustrated concepts only
 - In general, cannot visualise very high dimensional data: \Rightarrow can't directly observe overfitting/underfitting
- Main symptom of **underfitting**:
 - Poor performance even on the training data
- Main symptom of **overfitting**:
 - *Much* better performance on the training data than on independent test data
 - (Slightly better performance is to be expected)



NUI Galway
OÉ Gaillimh

Topic 2: Information-based Learning

Part 10: Review of topic



Learning Objectives Review

After completing this topic successfully, you will be able to ...

1. Explain what supervised learning is
2. Distinguish it from unsupervised learning and reinforcement learning
3. Describe in detail an algorithm for decision tree induction
4. Demonstrate the application of decision tree induction to a data set
5. List related algorithms
6. Discuss high-level concepts such as choice of hypothesis language, overfitting, underfitting and noise