

Lecture 4: Decision Tree Winter 2018

Kai-Wei Chang

CS @ UCLA

kw+cm146@kwchang.net

The instructor gratefully acknowledges Dan Roth, Sriram Sankararaman, Fei Sha, Ameet Talwalkar, Eric Eaton, and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

Key issues in machine learning

- ❖ Modeling
 - ❖ How to formulate your problem?
- ❖ Representation
 - ❖ What is the input/output space?
 - ❖ What is the hypothesis space?
- ❖ Algorithms
 - ❖ How to find the best hypothesis?

What Did We Learn?

- ❖ Learning problem:
Find a function that best separates the data

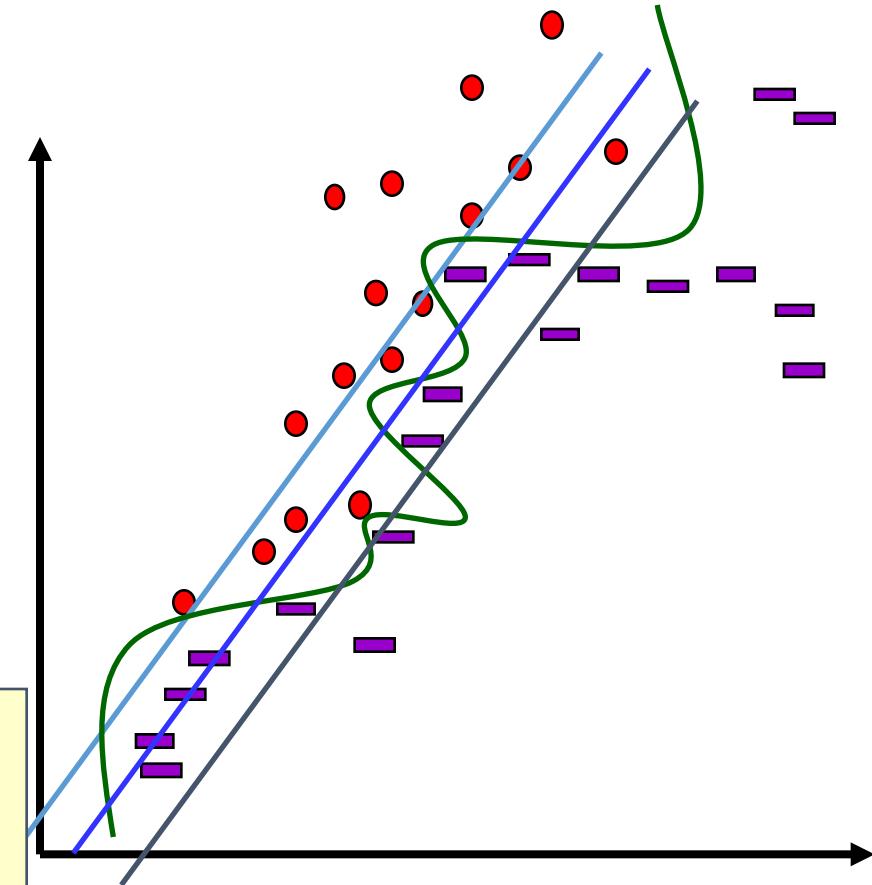
- ❖ What function?
- ❖ What's best?
- ❖ How to find it?

Linear:

x = data representation;

w = the classifier

$$Y = \text{sgn} \{w^T x\}$$



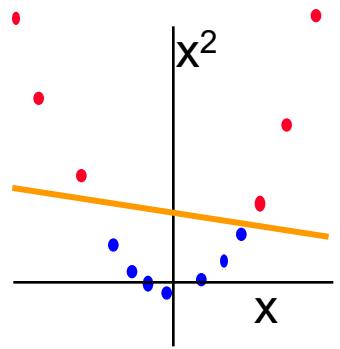
- ❖ A possibility: Define the learning problem to be:
Find a (linear) function that best separates the data

Motivation

- ❖ **Question 1:** Our solution learns a linear function; in principle, the target function may not be linear
- ❖ **Can we learn a function that is more flexible in terms of what it does with the feature space?**
- ❖ **Question 2:** Can we say something about the quality of what we learn (sample complexity, time complexity; quality)

Decision Trees

- ❖ Earlier, we decoupled the generation of the feature space from the learning.
 - ❖ Argued that we can map the examples into another space, in which the target functions are linearly separable.
 - ❖ How do we determine what are good mappings?
-
- ❖ The study of **decision trees** may shed some light on this.
 - ❖ Learning is done directly from the given data representation.
 - ❖ The algorithm ``transforms'' the data itself.

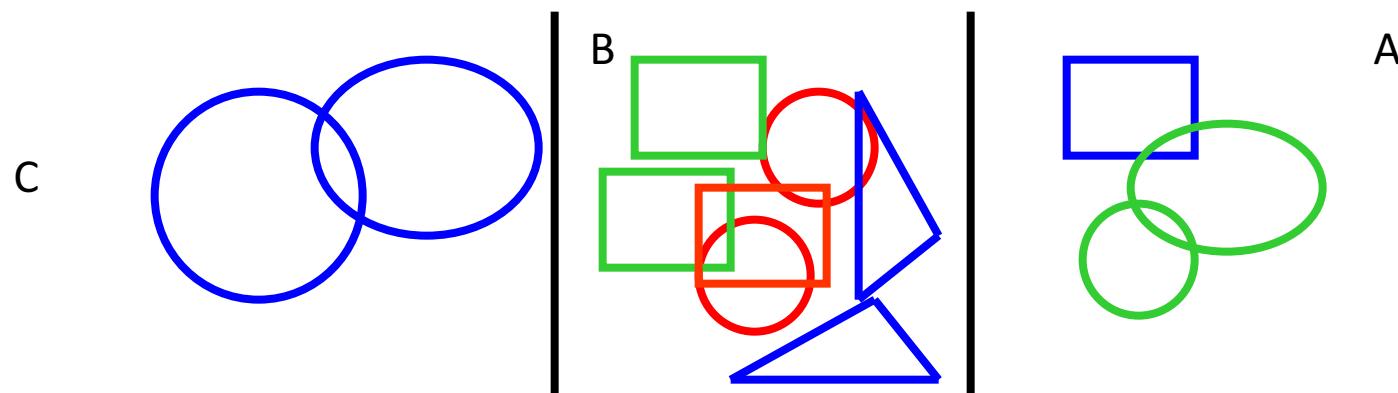


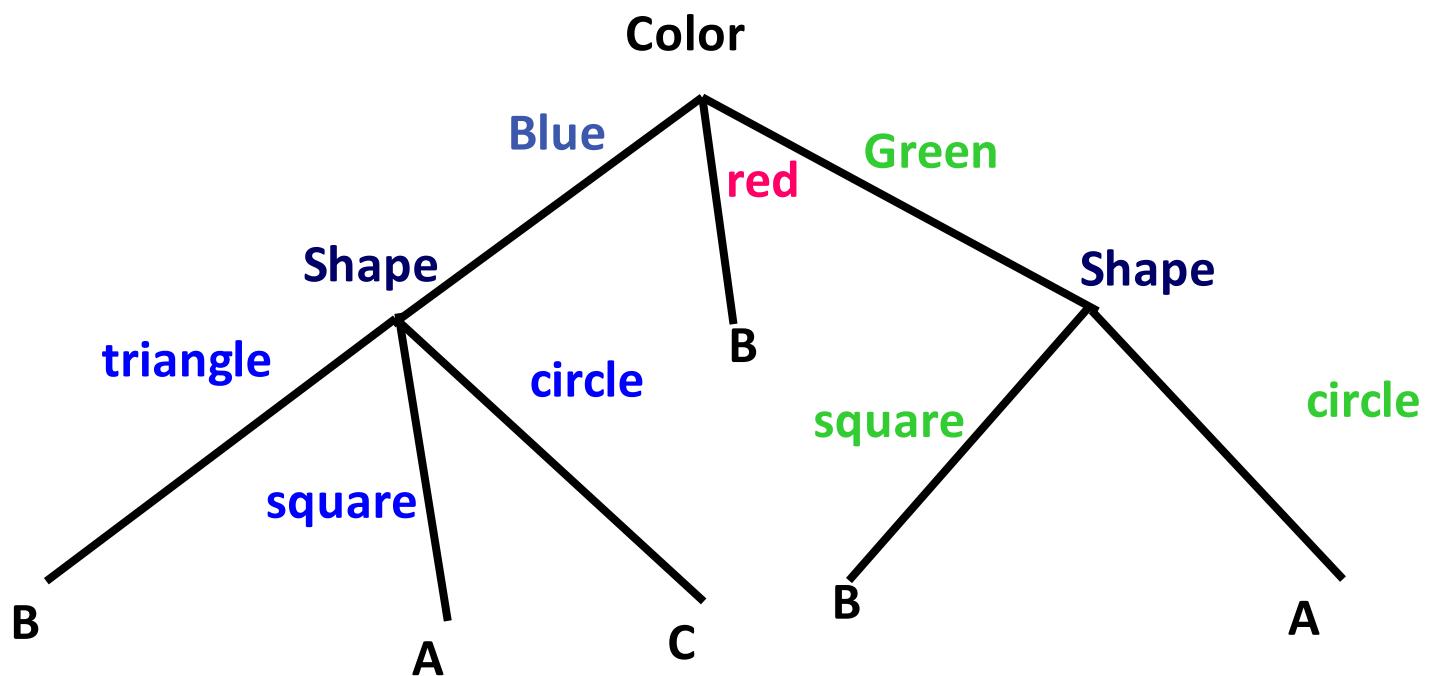
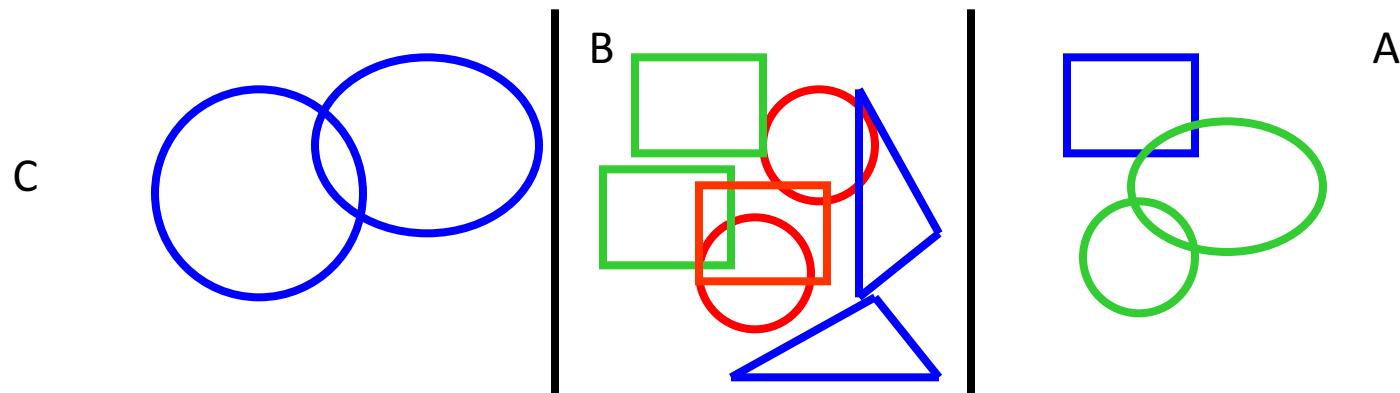
This Lecture

- ❖ Model/Representation: Decision trees
 - ❖ Non-linear classifiers
- ❖ Algorithm: Learning decision trees (ID3 algorithm)
 - ❖ Greedy heuristic (based on information gain)
Originally developed for discrete features
 - ❖ Some extensions to the basic algorithm

Sample dataset

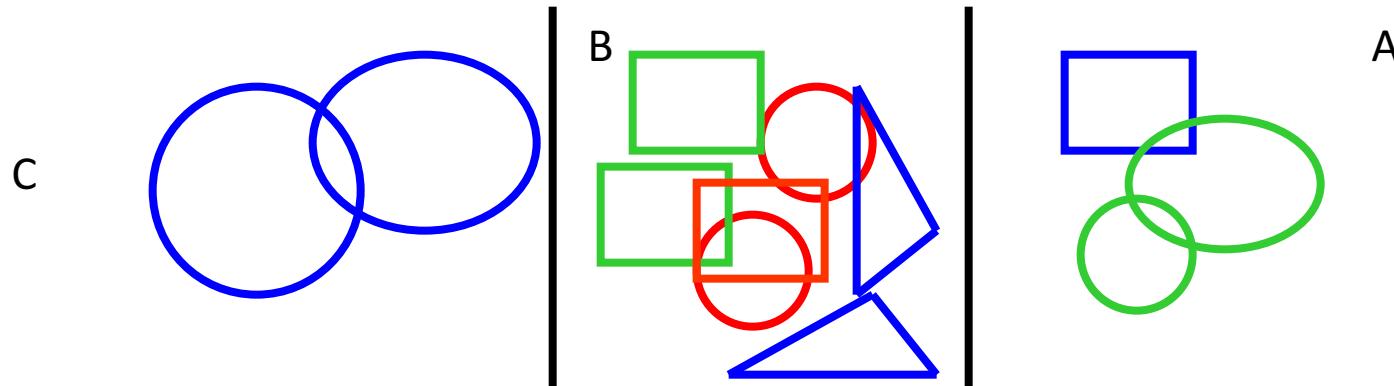
- ❖ What features we can used?
- ❖ What is the label for a **red triangle**?



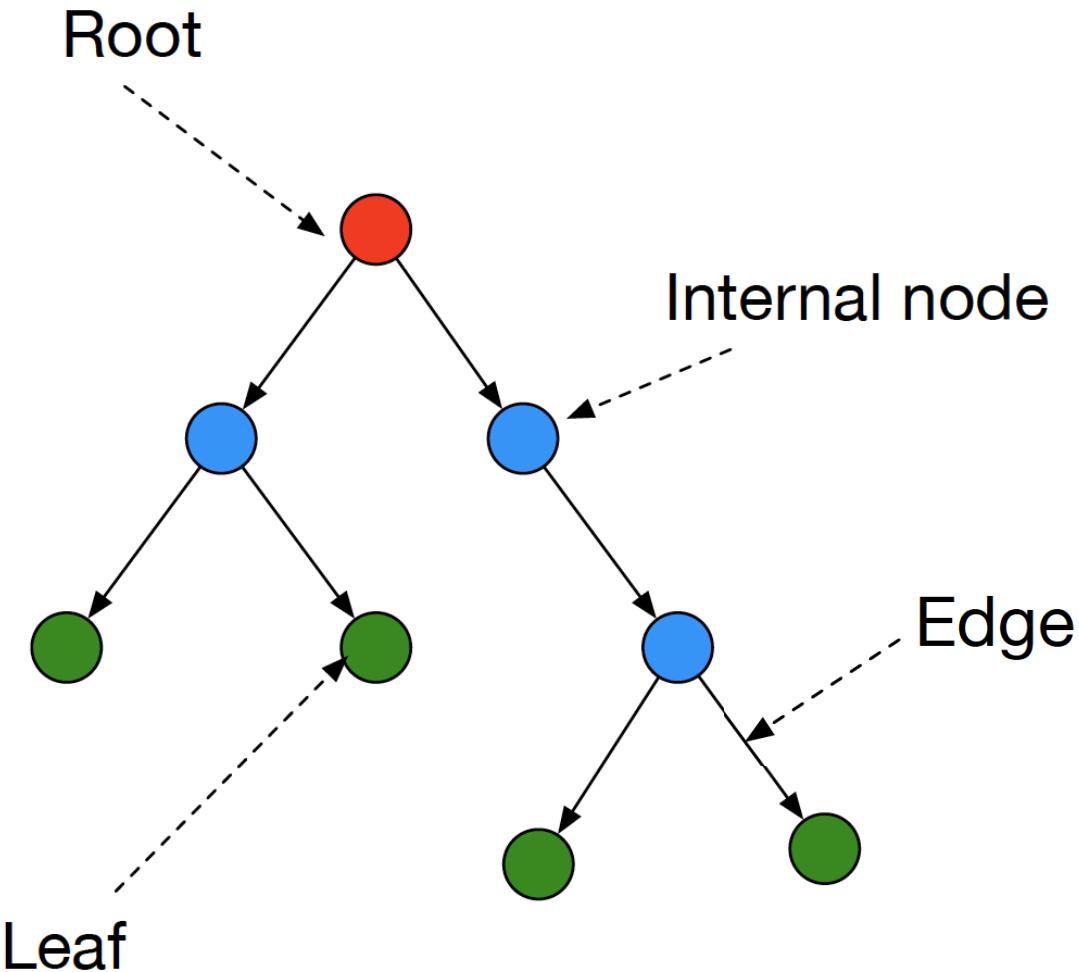


Decision Trees

- ❖ A hierarchical data structure that **represents data** by implementing a divide and conquer strategy
- ❖ Can be used as a classification or regression method
- ❖ Given a collection of examples, **learn a decision tree that represents it.**
- ❖ Use this representation to **classify new examples**



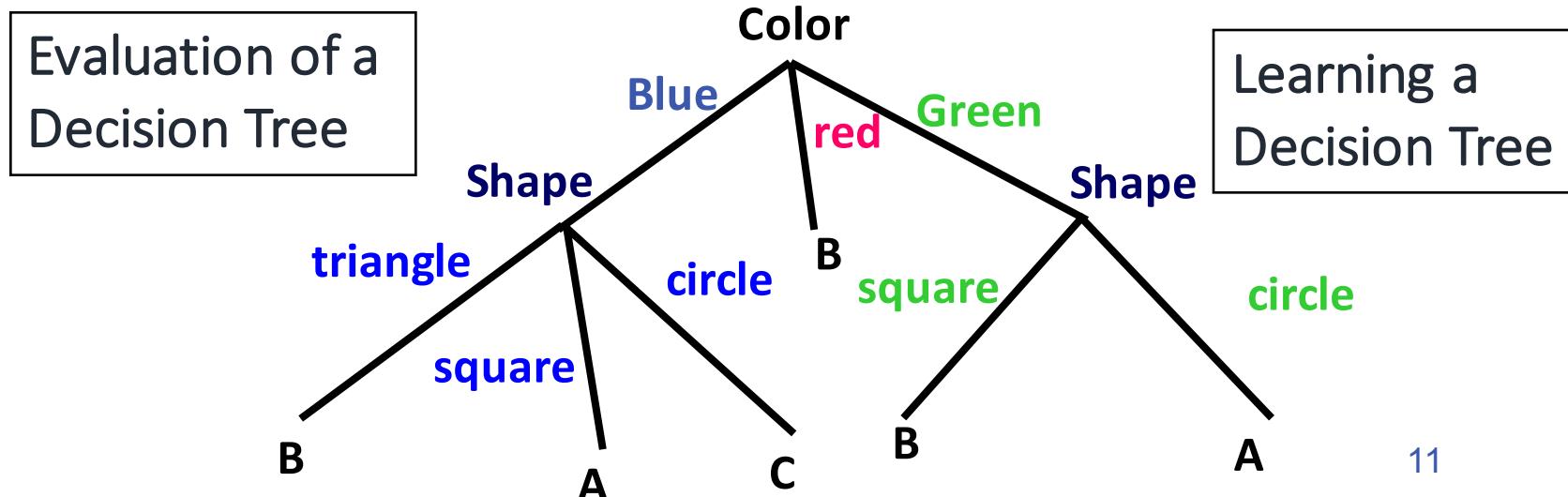
Terminology



Will sometimes drop the arrows on the edges

The Representation

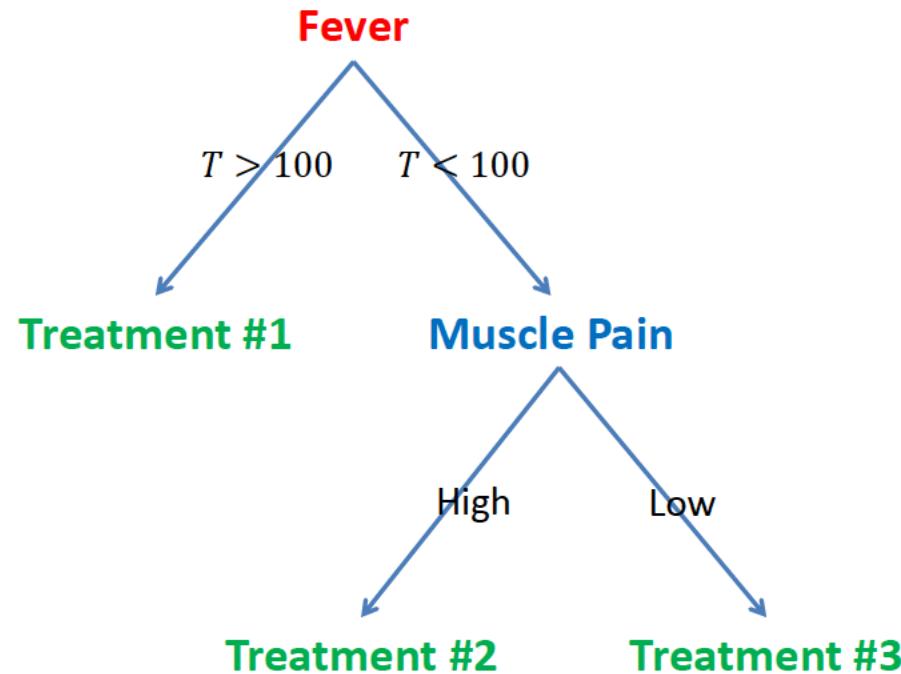
- ❖ Decision Trees are classifiers for instances represented as feature vectors (color= ; shape= ; label=)
- ❖ Nodes are **tests** for feature values
- ❖ Edges: There is one branch for each value of the feature
- ❖ Leaves specify the category (labels)
- ❖ Can categorize instances into multiple disjoint categories



Motivations:

Many decisions are tree structures

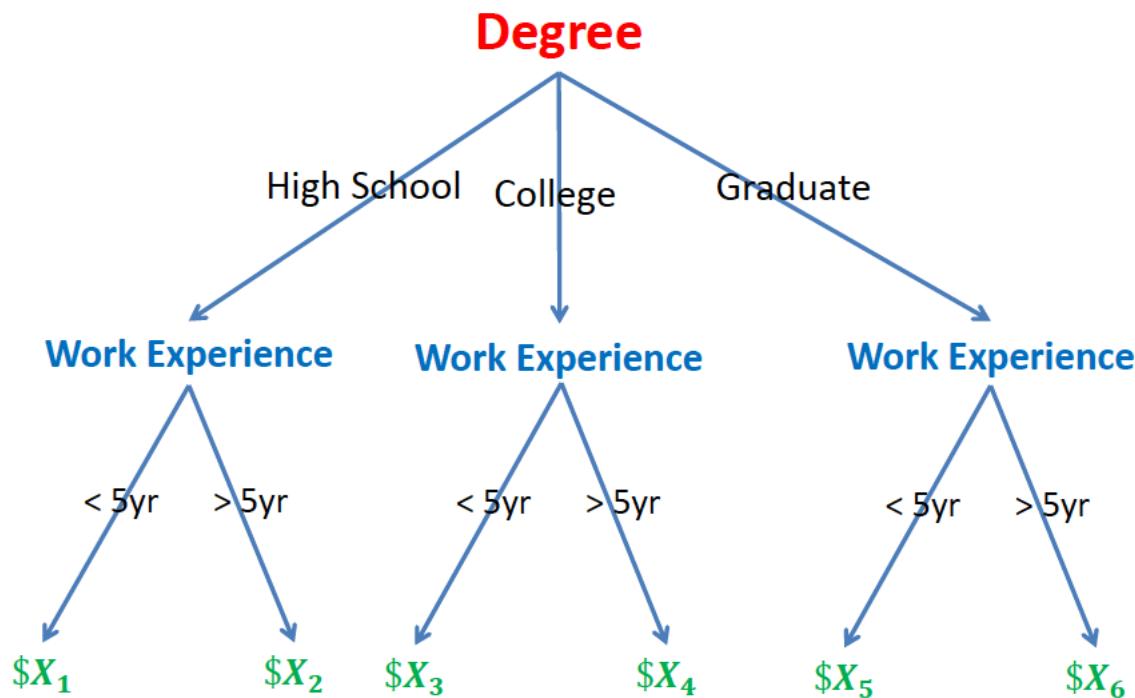
Medical treatment



Motivations:

Many decisions are tree structures

Salary in a company

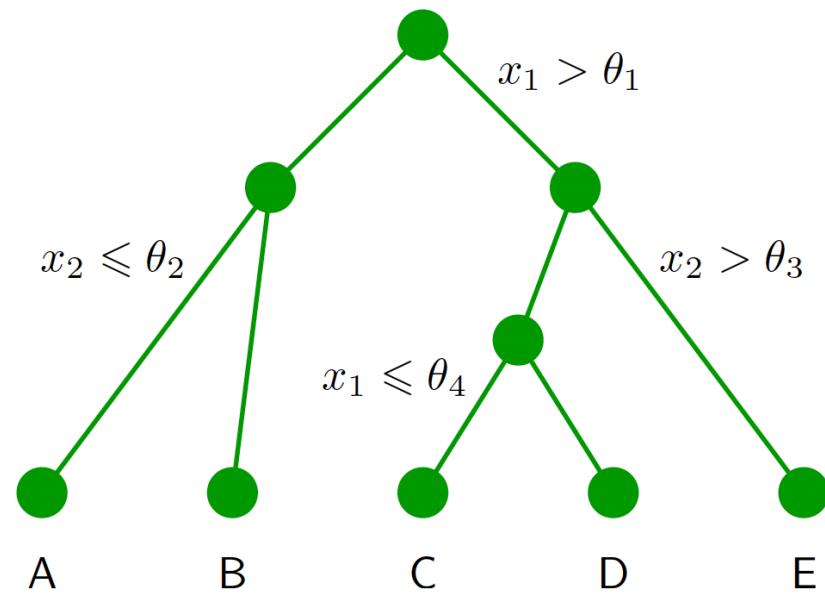
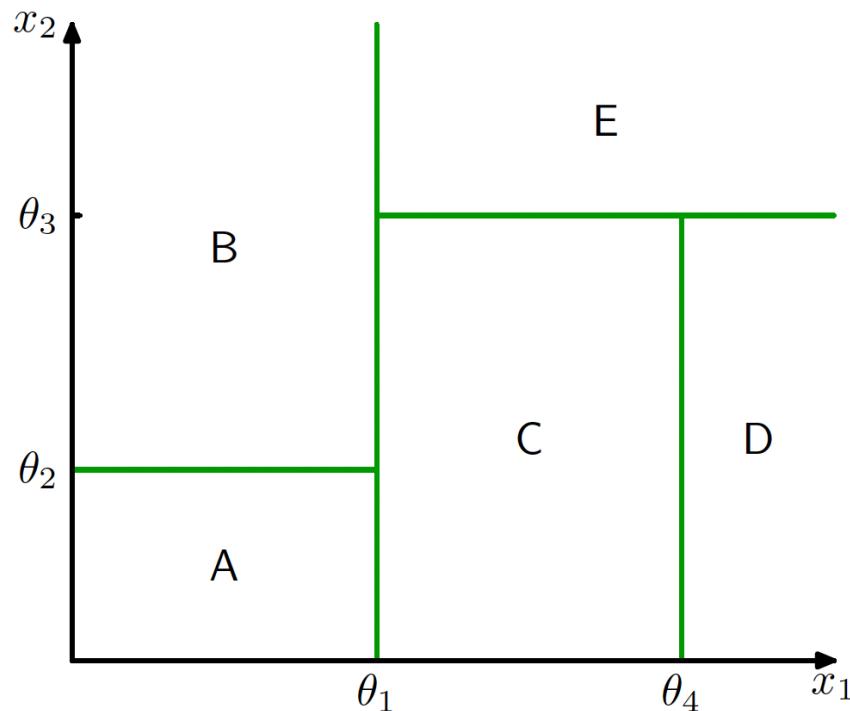


Decision Boundaries

-- Handling real-valued features

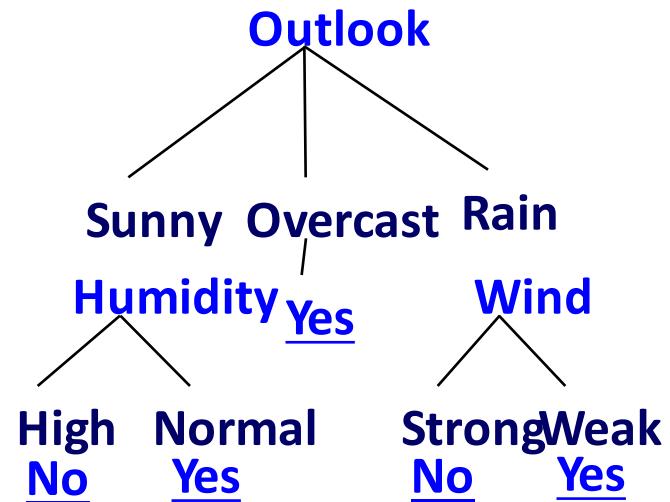
- ❖ Usually, instances are represented as attribute-value pairs (color=blue, shape = square, +)
- ❖ Numerical values can be used either by discretizing or by using thresholds for splitting nodes
- ❖ In this case, the tree divides the features space into axis-parallel rectangles, each labeled with one of the labels

A tree partitions the feature space



Advantages of Decision tree

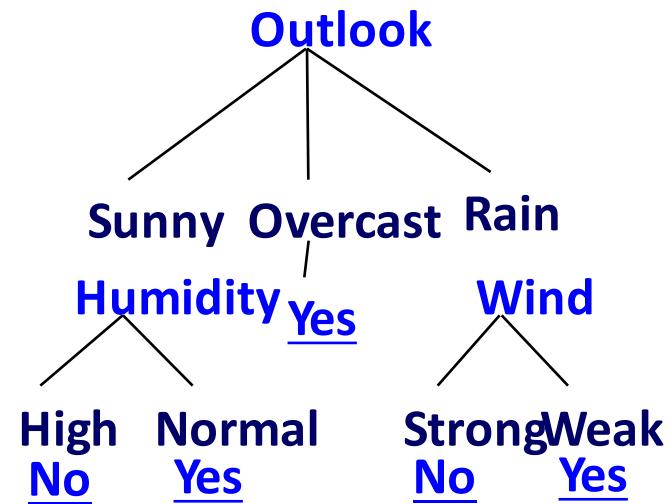
- ❖ Can represent any Boolean Function
- ❖ Can be viewed as a way to compactly represent a lot of data
- ❖ Natural representation
- ❖ The **evaluation** of the Decision Tree Classifier is easy



Challenge

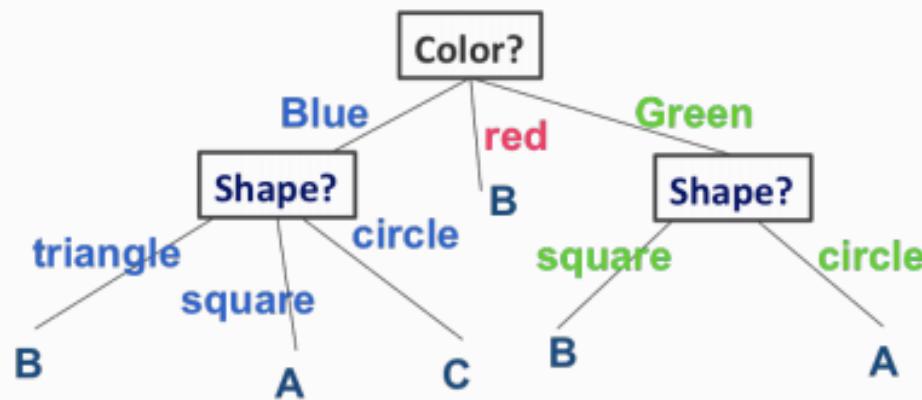
- ❖ Clearly, given data, there are many ways to represent it as a decision tree.
- ❖ Learning a **good** representation from data is the challenge.

A tree partitions the feature space



Expressivity of Decision Trees

- ❖ What Boolean functions can decision trees represent?
 - any Boolean function



(Color=blue AND Shape=triangle \Rightarrow Label=B) AND
(Color=blue AND Shape=square \Rightarrow Label=A) AND
(Color=blue AND Shape=circle \Rightarrow Label=C)....

Decision Trees

- ❖ Output is a discrete category. Real valued outputs are possible (regression trees)
- ❖ There are **efficient** algorithms for processing large amounts of data (but not too many features)
- ❖ There are methods for handling **noisy data** (classification noise and attribute noise) and for handling **missing attribute values**

Learning a decision tree

This Lecture

- ❖ Model/Representation: Decision trees
 - ❖ Non-linear classifiers
- ❖ Algorithm: Learning decision trees (ID3 algorithm)
 - ❖ Greedy heuristic (based on information gain)
Originally developed for discrete features
 - ❖ Some extensions to the basic algorithm

Will I play tennis today?

❖ Features

- ❖ Outlook: {Sun, Overcast, Rain}
- ❖ Temperature: {Hot, Mild, Cool}
- ❖ Humidity: {High, Normal, Low}
- ❖ Wind: {Strong, Weak}

❖ Labels

- ❖ Binary classification task: $Y = \{+, -\}$

Will I play tennis today?

	O	T	H	W	Play?
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

Outlook: S(unny),
O(vercast),
R(ainy)

Temperature: H(ot),
M(edium),
C(ool)

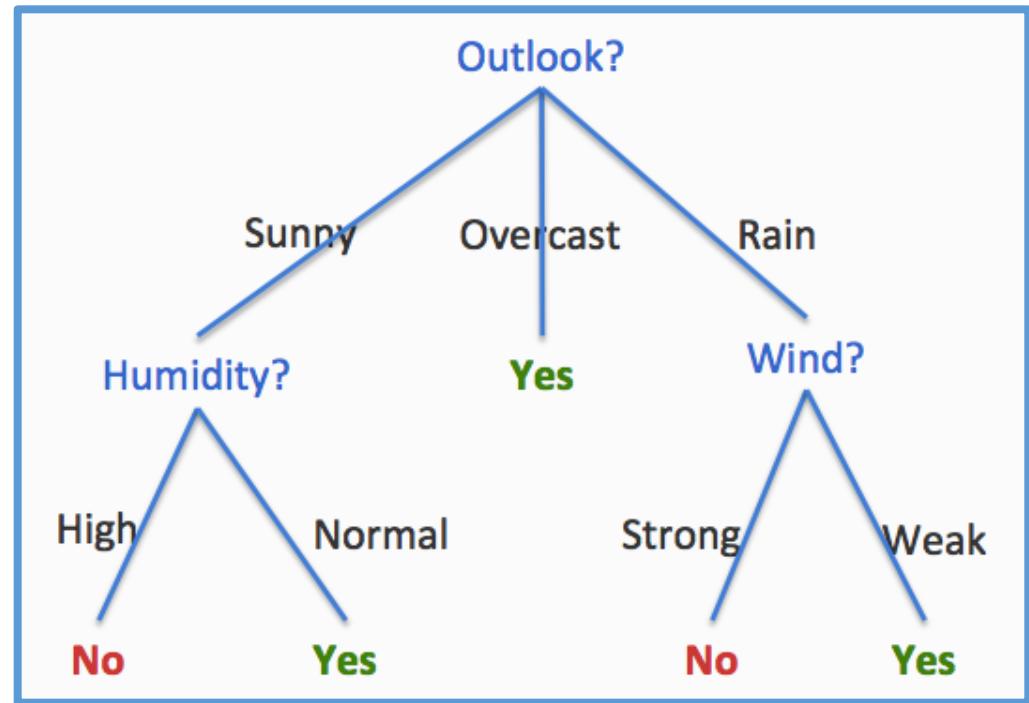
Humidity: H(igh),
N(ormal),
L(ow)

Wind: S(trong),
W(eak)

Basic Decision Trees Learning Algorithm

- ❖ Data is processed in Batch
(i.e. all the data available)
- ❖ Recursively build a decision tree top down.

O	T	H	W	Play?	
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

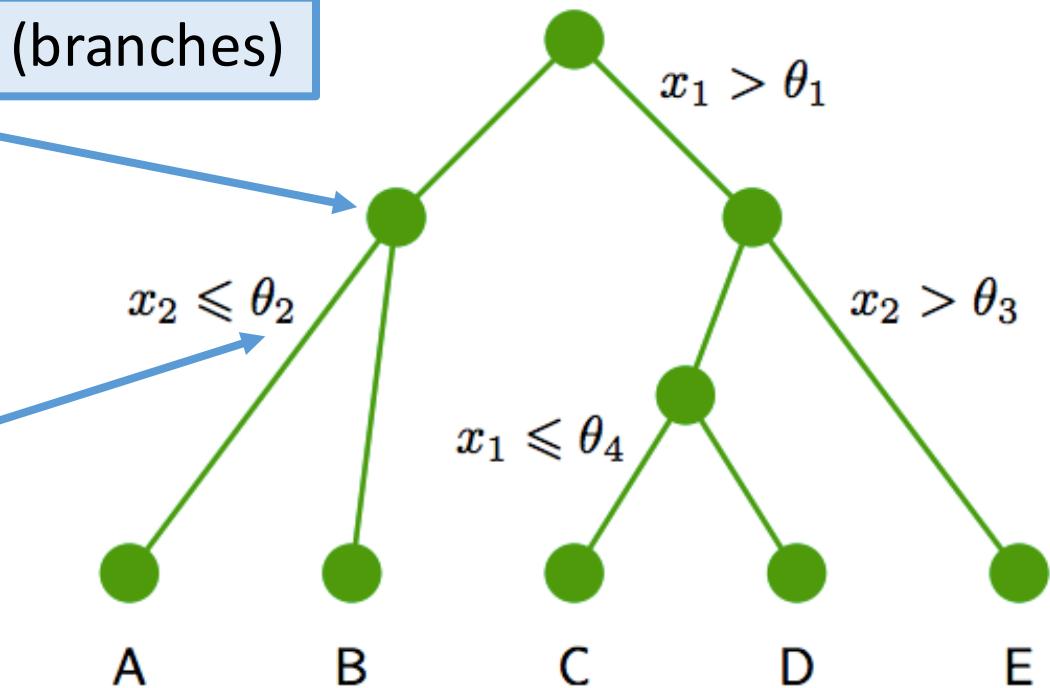


What do we need to learn?

The structure of the tree (branches)

The threshold values

Values for the leaves



What do we need to learn?

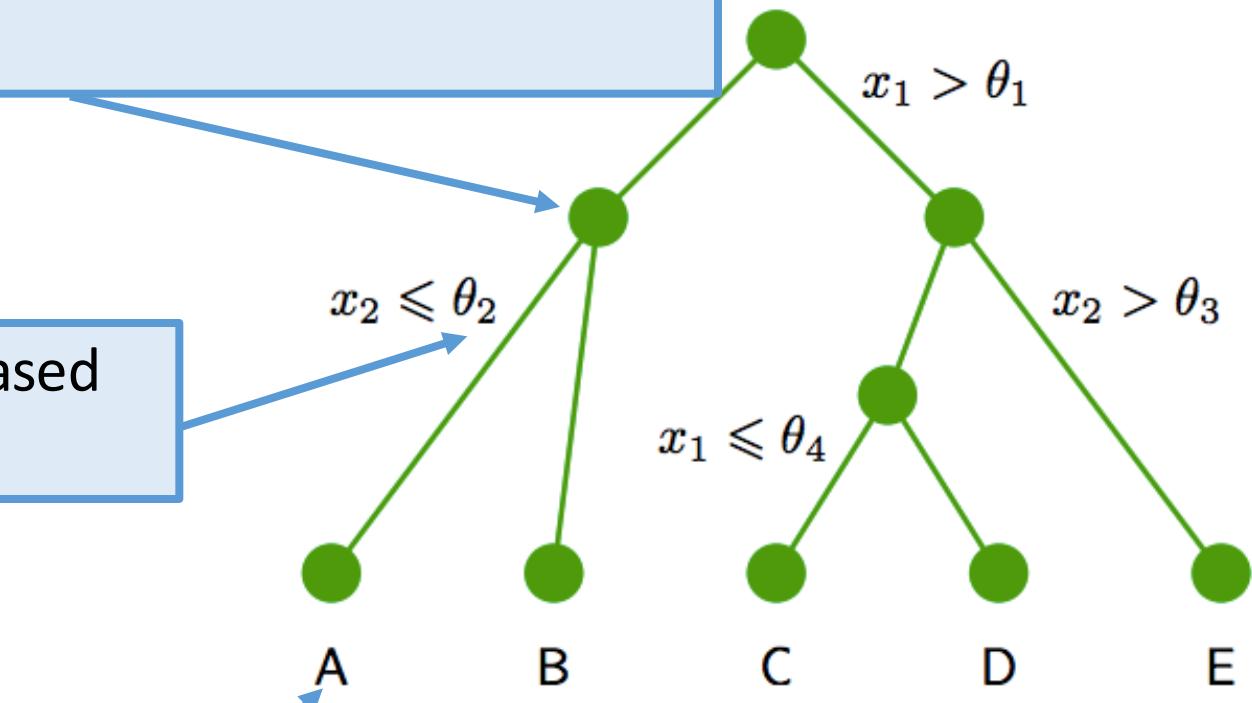
Pick one feature that **best** classifier the data

How?

Create branches based
on feature values

Values for the leaves =
the output label.

When to stop?



DT algorithm: ID3(S , Attributes, Label)

- ❖ A recursive algorithm
- ❖ Recursively build a decision tree top down.
- ❖ Base case:
 - If all examples are labeled the same
Return a single node tree with Label
 - Otherwise
 - Recursive decision tree algorithm
(see next slide)

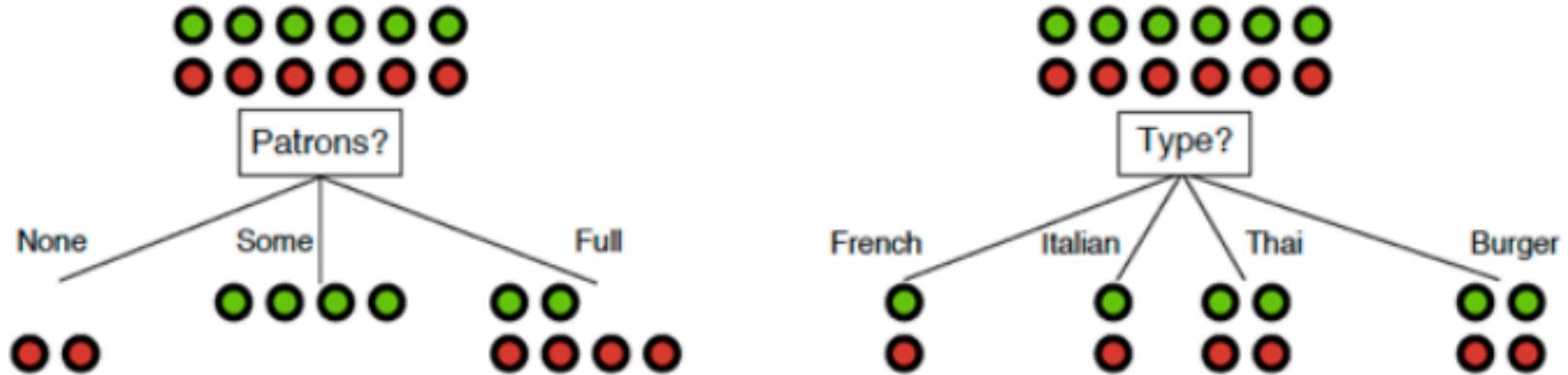
DT algorithm: ID3(S , Attributes, Label)

1. If all examples have a same label
return a single node tree with LabelBase case
2. Create a Root node for tree
3. $A = \text{attribute in Attributes that } \underline{\text{best}} \text{ classifies } S$
4. For each possible value v of A
 1. Add a new tree branch corresponding to $A=v$
 2. Let S_v be the subset of examples in S with $A=v$
 3. if S_v is empty:
add leaf node with the common value of Label in S
Else: below this branch add the subtree
 $\text{ID3}(S_v, \text{Attributes} - \{a\}, \text{Label})$
4. Return Root

Which attribute to split?

- ❖ The goal is to have the resulting decision tree as small as possible
- ❖ But, finding the minimal decision tree consistent with the data is NP-hard
- ❖ The recursive algorithm is a greedy heuristic search for a simple tree, but cannot guarantee optimality.
- ❖ The main decision in the algorithm is the selection of the next attribute to condition on.

Which attribute to split?



Patrons? is a better choice—gives **information** about the classification

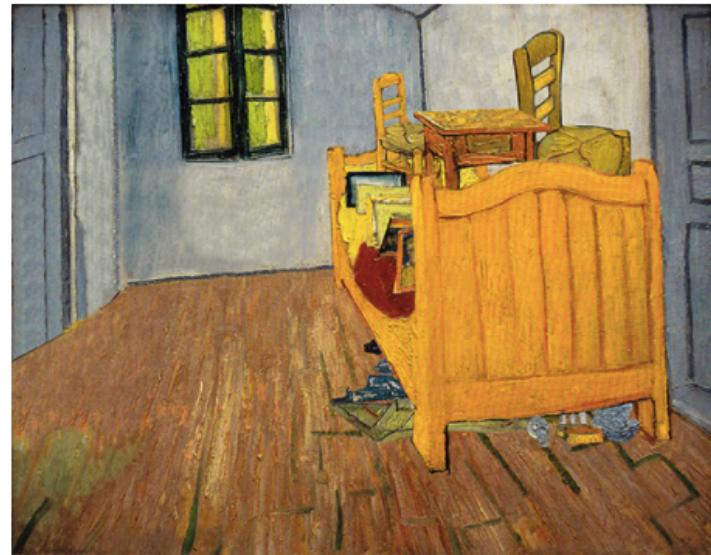
Need a way to quantify things

Which attribute to split?

- ❖ The goal is to have the resulting decision tree as small as possible
- ❖ The main decision in the algorithm is the selection of the next attribute to condition on.
- ❖ We want attributes that split the examples to sets that are **relatively pure in one label**; this way we are closer to a leaf node.
- ❖ The most popular heuristics is based on **information gain**, originated with the ID3 system of Quinlan.

How to measure information gain?

- ❖ Idea: Gaining information reduces uncertainty
- ❖ Uncertainty can be measured by Entropy



Vincent Van Gogh: Bedroom in Arles

By Ursus Wehrli

High entropy

Low entropy

How to measure information gain?

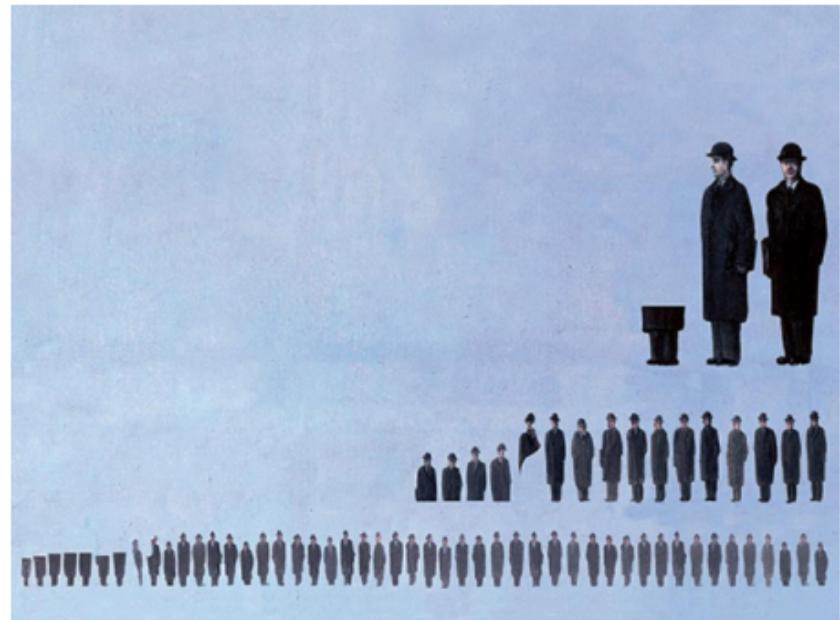
- ❖ Idea: Gaining information reduces uncertainty
- ❖ Uncertainty can be measured by Entropy

René Magritte "Golconda"



High entropy

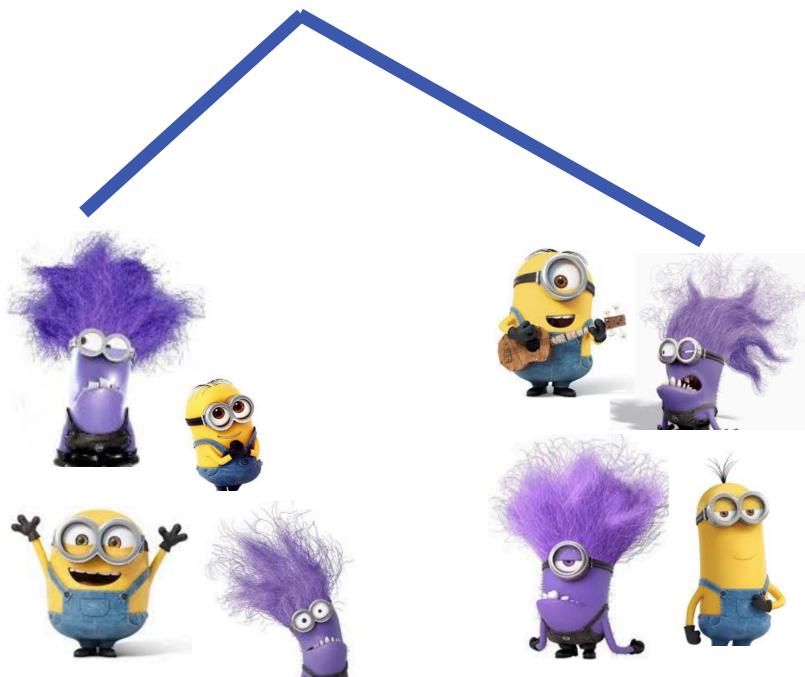
By Ursus Wehrli



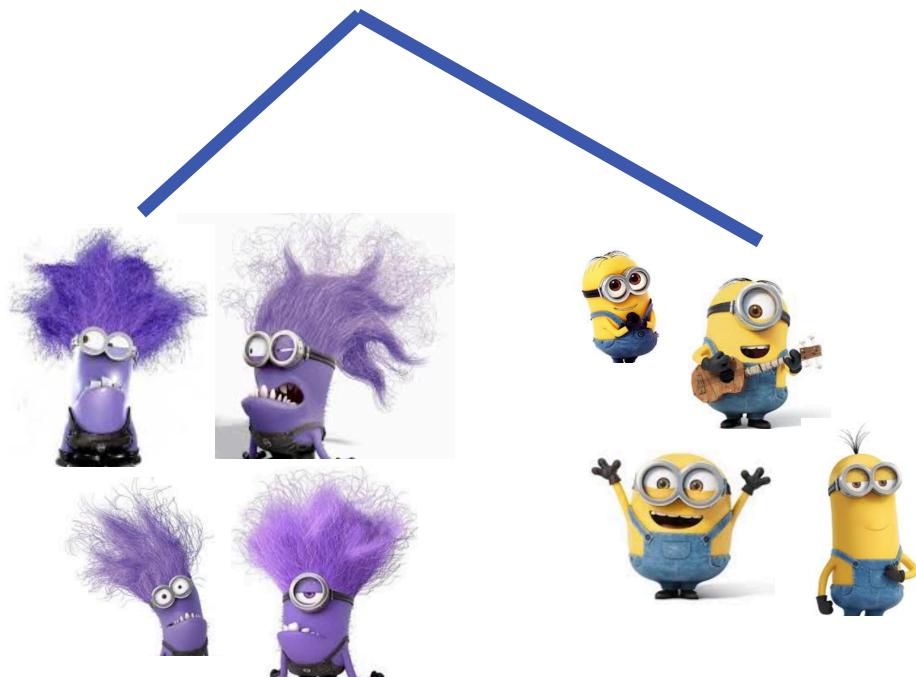
Low entropy

How to measure information gain?

- ❖ Idea: Gaining information reduces uncertainty
- ❖ Uncertainty can be measured by Entropy



High entropy



Low entropy

Entropy

- ❖ Entropy (impurity, disorder) of a set of examples, S , relative to a binary classification is:
$$H[S] = -P_+ \log_2(P_+) - P_- \log_2(P_-)$$
- ❖ where P_+ is the proportion of positive examples in S and P_- is the proportion of negatives.

Here we define $0 \log 0 = 0$

Entropy (formal definition)

- ❖ If a random variable S has K different values, a_1, a_2, \dots, a_K , its entropy is given by

$$H[S] = - \sum_{v=1}^K P(S = a_v) \log_2 P(S = a_v)$$

- ❖ Measures the amount of uncertainty of a random variable with a specific probability distribution. Higher it is, less confident we are in its outcome

Entropy (intuition)

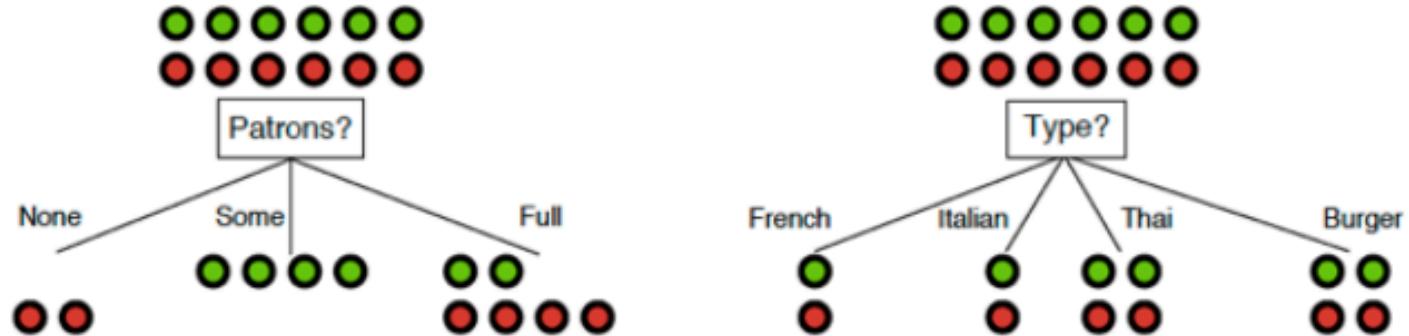
- ❖ In average, how many bits do we need to send the message (#bits/#length of message)



Entropy (intuition)

- ❖ In average, how many bits do we need to send the message (#bits/#length of message)
- ❖ Consider you have four possible tokens (a,b,c,d). What is the best way to encode them?
- ❖ All examples belong to the same category
 - e.g., aaaaaaaaaaaaaaaaaaaaaaaaaaaaa
 - no need to communicate
- ❖ If all the examples are equally mixed (0.5, 0.5):
 - e.g., abbacaccddd.....
 - two bits for each token: (a:00, b:01, c:10, d:11)
- ❖ If $\frac{1}{4}$ of message is a, and $\frac{1}{2}$ is b and $\frac{1}{4}$ is c in average:
 - e.g., abbbbacc.....
 - (a:00, b:1, c:01, d:--)

Which attribute to split



Patrons? is a better choice—gives **information** about the classification

Patron vs. Type?

By choosing Patron, we end up with a partition (3 branches) with smaller entropy, ie, smaller uncertainty (0.45 bit)

By choosing Type, we end up with uncertainty of 1 bit.

Thus, we choose Patron over Type.

Uncertainty if we go with “Patron”

For “None” branch

$$-\left(\frac{0}{0+2} \log \frac{0}{0+2} + \frac{2}{0+2} \log \frac{2}{0+2}\right) = 0$$

For “Some” branch

$$-\left(\frac{4}{4+0} \log \frac{4}{4+0} + \frac{4}{4+0} \log \frac{4}{4+0}\right) = 0$$

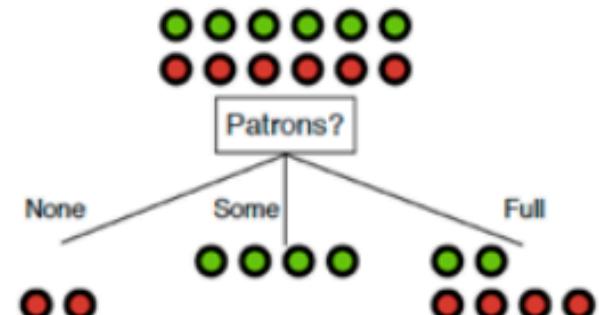
For “Full” branch

$$-\left(\frac{2}{2+4} \log \frac{2}{2+4} + \frac{4}{2+4} \log \frac{4}{2+4}\right) \approx 0.9$$

For choosing “Patrons”

weighted average of each branch: this quantity is called **conditional entropy**

$$\frac{2}{12} * 0 + \frac{4}{12} * 0 + \frac{6}{12} * 0.9 = 0.45$$



Information Gain

- ❖ The information gain of an attribute a is the expected reduction in entropy caused by partitioning on this attribute

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

- ❖ S_v is the subset of S for which attribute a has value v .
- ❖ The entropy of partitioning the data is calculated by weighing the entropy of each partition by its size relative to the original set

Will I play tennis today?

	O	T	H	W	Play?
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

Outlook: S(unny),
O(vercast),
R(ainy)

Temperature: H(ot),
M(edium),
C(ool)

Humidity: H(igh),
N(ormal),
L(ow)

Wind: S(strong),
W(eak)

Will I play tennis today?

	O	T	H	W	Play?
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

Current entropy:

$$p = 9/14$$

$$n = 5/14$$

$$H(\text{Play?}) = -(9/14) \log_2(9/14)$$

$$-(5/14) \log_2(5/14)$$

$$= 0/.94$$

Information Gain: Outlook

	O	T	H	W	Play?
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

Information Gain: Outlook

	O	T	H	W	Play?
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

Outlook = sunny: 5 of 14 examples

$$p = 2/5 \quad n = 3/5 \quad H_s = 0.971$$

Information Gain: Outlook

	O	T	H	W	Play?
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

Outlook = sunny: 5 of 14 examples

$$p = 2/5 \quad n = 3/5 \quad H_s = 0.971$$

Outlook = overcast: 4 of 14 examples

$$p = 4/4 \quad n = 0 \quad H_o = 0$$

Information Gain: Outlook

	O	T	H	W	Play?
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

Outlook = sunny: 5 of 14 examples

$$p = 2/5 \quad n = 3/5 \quad H_s = 0.971$$

Outlook = overcast: 4 of 14 examples

$$p = 4/4 \quad n = 0 \quad H_o = 0$$

Outlook = rainy: 5 of 14 examples

$$p = 3/5 \quad n = 2/5 \quad H_R = 0.971$$

Expected entropy:

$$(5/14) \times 0.971 + (4/14) \times 0 + (5/14) \times 0.971 \\ = 0.694$$

Information gain:

$$0.940 - 0.694 = 0.246$$

Information Gain: Humidity

O	T	H	W	Play?
1	S	H	W	-
2	S	H	S	-
3	O	H	W	+
4	R	M	W	+
5	R	C	W	+
6	R	C	S	-
7	O	C	S	+
8	S	M	W	-
9	S	C	W	+
10	R	M	W	+
11	S	M	S	+
12	O	M	S	+
13	O	H	W	+
14	R	M	S	-

Information Gain: Humidity

O	T	H	W	Play?
1	S	H	W	-
2	S	H	S	-
3	O	H	W	+
4	R	M	W	+
5	R	C	W	+
6	R	C	S	-
7	O	C	S	+
8	S	M	W	-
9	S	C	W	+
10	R	M	W	+
11	S	M	S	+
12	O	M	S	+
13	O	H	W	+
14	R	M	S	-

Humidity = High:

$$p = 3/7 \quad n = 4/7$$

$$H_h = 0.985$$

Information Gain: Humidity

O	T	H	W	Play?
1	S	H	W	-
2	S	H	S	-
3	O	H	W	+
4	R	M	W	+
5	R	C	W	+
6	R	C	S	-
7	O	C	S	+
8	S	M	W	-
9	S	C	W	+
10	R	M	W	+
11	S	M	S	+
12	O	M	S	+
13	O	H	W	+
14	R	M	S	-

Humidity = High:

$$p = 3/7 \quad n = 4/7 \quad H_h = 0.985$$

Humidity = Normal:

$$p = 6/7 \quad n = 1/7 \quad H_o = 0.592$$

Expected entropy:

$$(7/14) \times 0.985 + (7/14) \times 0.592 = 0.7885$$

Information Gain: Humidity

O	T	H	W	Play?
1	S	H	W	-
2	S	H	S	-
3	O	H	W	+
4	R	M	W	+
5	R	C	W	+
6	R	C	S	-
7	O	C	S	+
8	S	M	W	-
9	S	C	W	+
10	R	M	W	+
11	S	M	S	+
12	O	M	S	+
13	O	H	W	+
14	R	M	S	-

Humidity = High:

$$p = 3/7 \quad n = 4/7 \quad H_h = 0.985$$

Humidity = Normal:

$$p = 6/7 \quad n = 1/7 \quad H_o = 0.592$$

Expected entropy:

$$(7/14) \times 0.985 + (7/14) \times 0.592 = 0.7885$$

Information gain:

$$0.940 - 0.7885 = 0.1515$$

Which feature to split on?

	O	T	H	W	Play?
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

Information gain:

Outlook: 0.246

Humidity: 0.151

Wind: 0.048

Temperature: 0.029

Which feature to split on?

	O	T	H	W	Play?
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

Information gain:

Outlook: 0.246

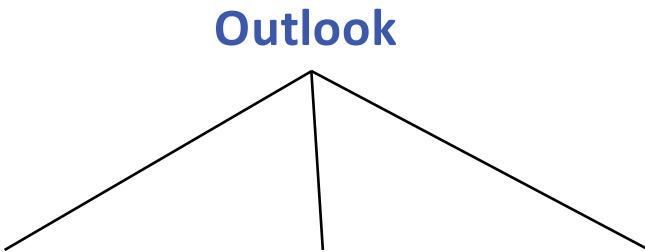
Humidity: 0.151

Wind: 0.048

Temperature: 0.029

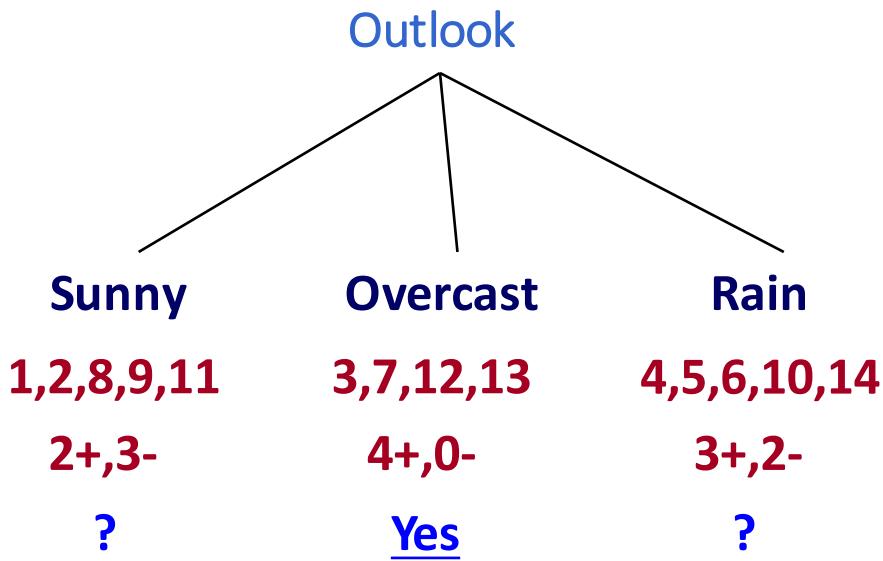
→ Split on Outlook

An Illustrative Example



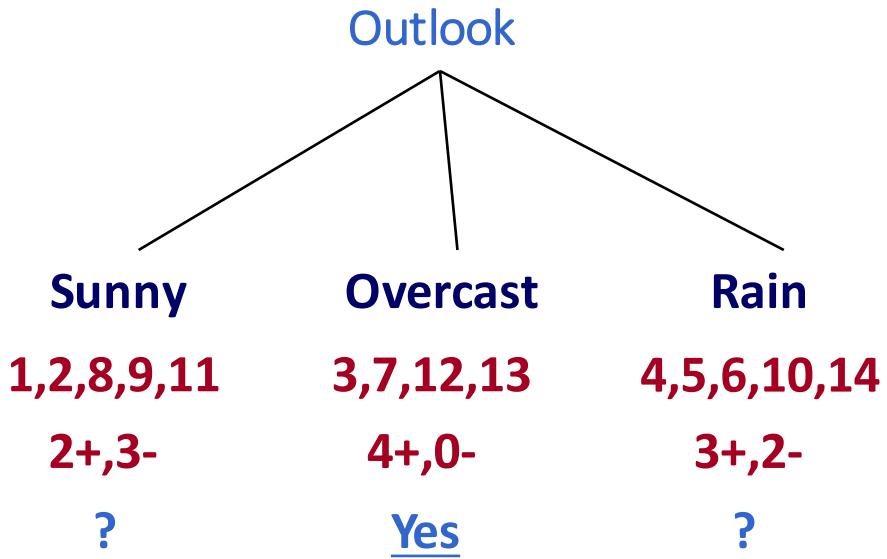
Gain(S, Humidity) = 0.151
Gain(S, Wind) = 0.048
Gain(S, Temperature) = 0.029
Gain(S, Outlook) = 0.246

An Illustrative Example



	O	T	H	W	Play?
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

An Illustrative Example

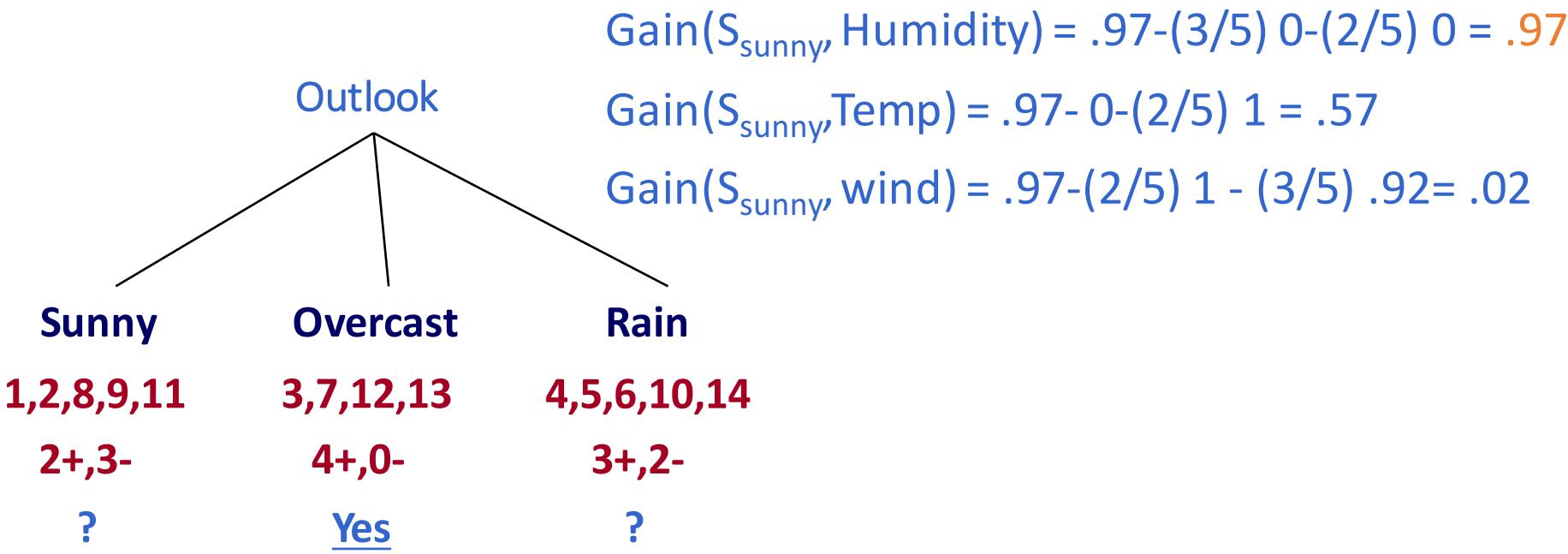


	O	T	H	W	Play?
1	S	H	H	W	-
2	S	H	H	S	-
3	O	H	H	W	+
4	R	M	H	W	+
5	R	C	N	W	+
6	R	C	N	S	-
7	O	C	N	S	+
8	S	M	H	W	-
9	S	C	N	W	+
10	R	M	N	W	+
11	S	M	N	S	+
12	O	M	H	S	+
13	O	H	N	W	+
14	R	M	H	S	-

Continue until:

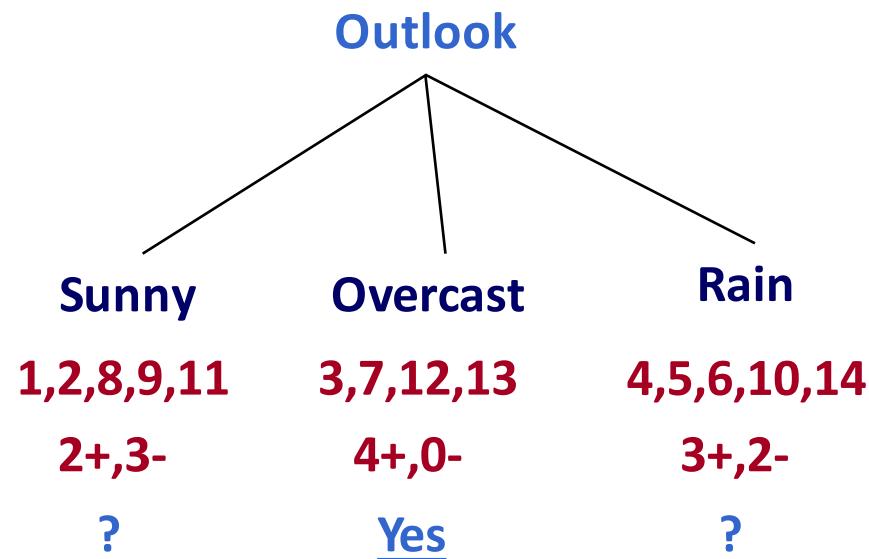
- Every attribute is included in path, or,
- All examples in the leaf have same label

An Illustrative Example

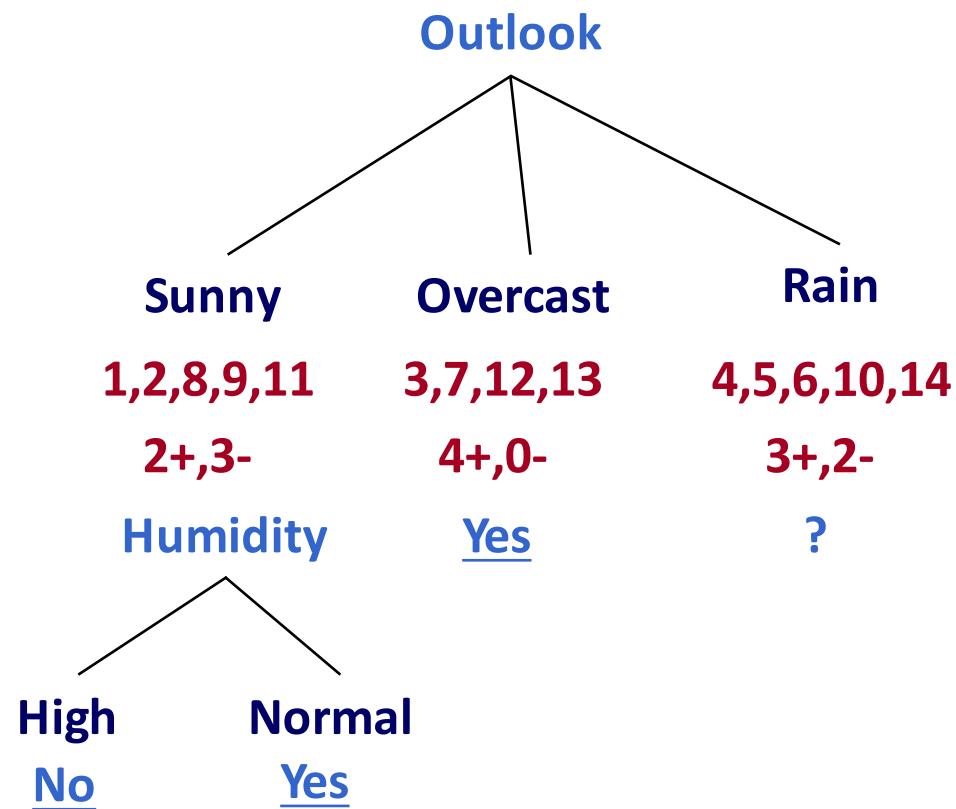


Day	Outlook	Temperature	Humidity	Wind	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes

An Illustrative Example



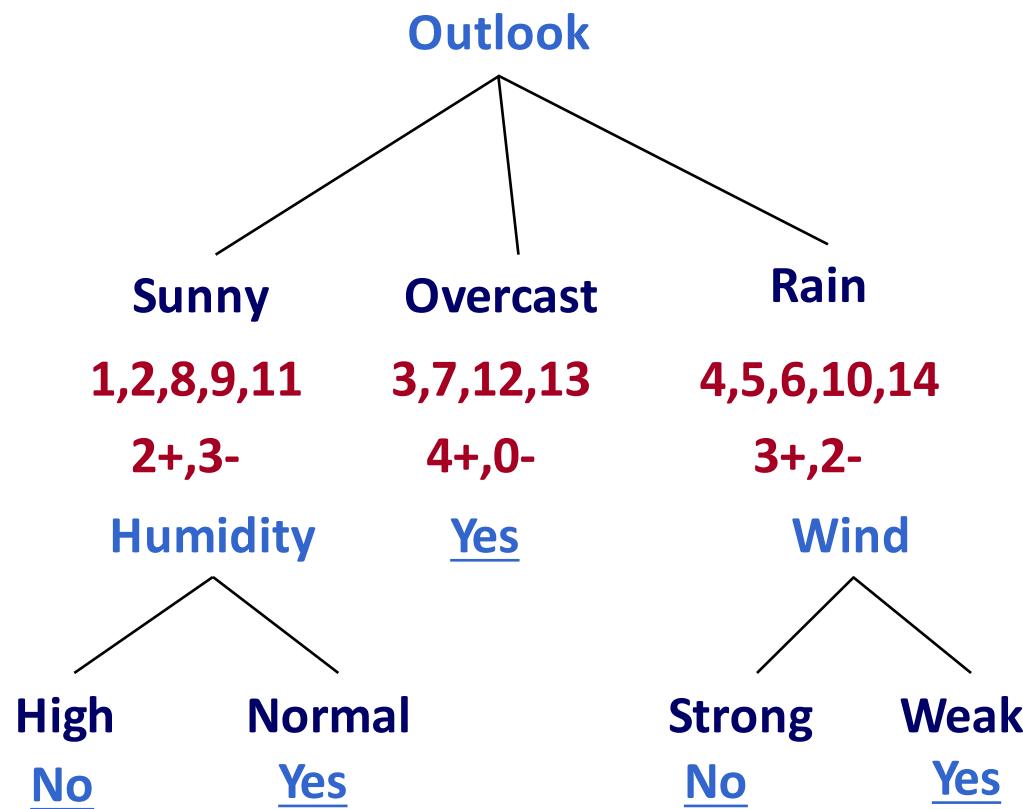
An Illustrative Example



induceDecisionTree(S)

1. Does S uniquely define a class?
if all $s \in S$ have the same label y : **return** S ;
2. Find the feature with the most information gain:
 $i = \operatorname{argmax}_i Gain(S, X_i)$
3. Add children to S :
for k **in** $\text{Values}(X_i)$:
 $S_k = \{s \in S \mid x_i = k\}$
 addChild(S , S_k)
 $\text{induceDecisionTree}(S_k)$
return S ;

An Illustrative Example



Summary: Learning Decision Trees

1. **Representation**: What are decision trees?

- ❖ A hierarchical data structure that represents data

2. **Algorithm**: Learning decision trees

The ID3 algorithm: A greedy heuristic

- ❖ If all the examples have the same label, create a leaf with that label
- ❖ Otherwise, find the “most informative” attribute and split the data for different values of that attribute
- ❖ Recurse on the splits