

# Decision Trees

Sven Koenig, USC

Russell and Norvig, 3<sup>rd</sup> Edition, Section 18.3

These slides are new and can contain mistakes and typos.  
Please report them to Sven ([skoenig@usc.edu](mailto:skoenig@usc.edu)).

## Rule Learning

- So far, we assumed that rules need to be specified by experts.
- Now, we find out how a system can learn them from examples.
- Thus, we study how to acquire knowledge with machine learning.

## Inductive Learning for Classification

- Labeled examples

How old are they?	What is their current salary per year?	Do they have a savings account?	Have they ever declared bankruptcy?	...	Would you issue a credit card to them?
52	\$150,000	yes	no	...	yes
40	\$50,000	no	yes	...	no
20	\$60,000	yes	no	...	yes
31	\$20,000	yes	no	...	yes

- Unlabeled examples

How old are they?	What is their current salary per year?	Do they have a savings account?	Have they ever declared bankruptcy?	...	Would you issue a credit card to them?
26	\$40,000	no	no	...	?

## Inductive Learning for Classification

- Labeled examples

Feature_1	Feature_2	Class
true	true	true
true	false	false
false	true	false

- Unlabeled examples

Feature_1	Feature_2	Class
false	false	?

## Inductive Learning for Classification

- Labeled examples

Feature_1	Feature_2	Class
true	true	true
true	false	false
false	true	false

Learn  $f(\text{Feature\_1}, \text{Feature\_2}) = \text{Class}$  from

$$f(\text{true}, \text{true}) = \text{true}$$

$$f(\text{true}, \text{false}) = \text{false}$$

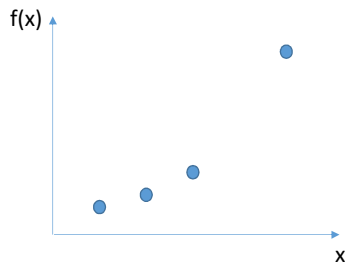
$$f(\text{false}, \text{true}) = \text{false}$$

The function needs to be consistent with all labeled examples and should make the fewest mistakes on the unlabeled examples.

- Unlabeled examples

Feature_1	Feature_2	Class
false	false	?

## Inductive Learning for Classification



- Labeled examples

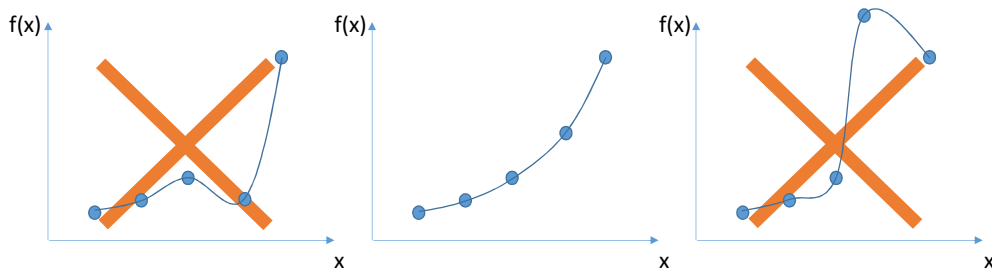
Feature_1 = x	Class = f(x)
1.0	0.5
2.0	0.7
3.0	1.0
5.0	3.0

- Unlabeled examples

Feature_1 = x	Class = f(x)
4.0	?

## Inductive Learning for Classification

- Function learning needs bias, i.e. to prefer some functions over others.

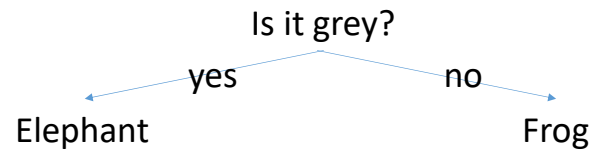


- Many students choose the function in the center.
- They prefer “simple” functions.

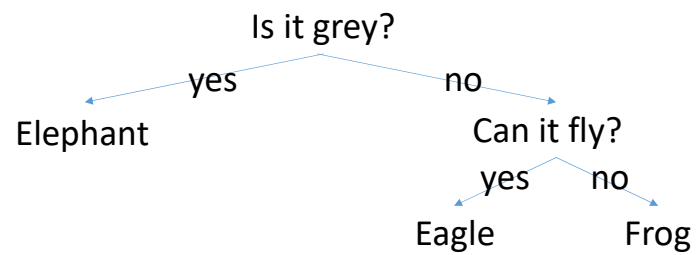
## Example: Decision Tree (and Rule) Learning

Frog

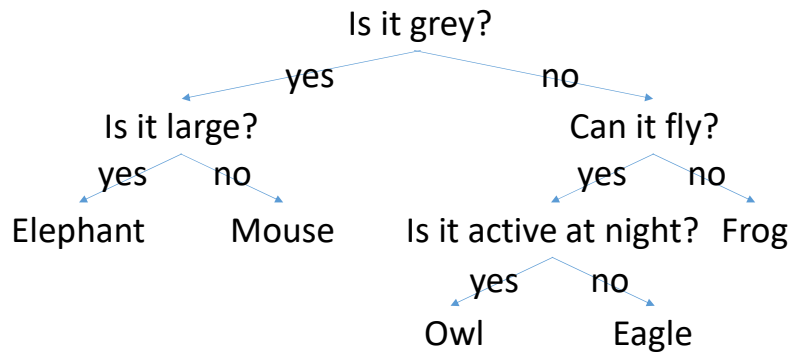
## Example: Decision Tree (and Rule) Learning



## Example: Decision Tree (and Rule) Learning



## Example: Decision Tree (and Rule) Learning

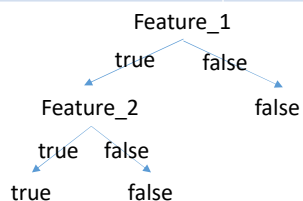


- Objective: Learn a decision tree
- From now on: binary (feature and class) values only.

## Example: Decision Tree (and Rule) Learning

- Labeled examples

Feature_1	Feature_2	Class
true	true	true
true	false	false
false	true	false



Feature\_1 AND Feature\_2 → Class  
 Feature\_1 AND NOT Feature\_2 → NOT Class  
 NOT Feature\_1 → NOT Class

- Unlabeled examples (note: classification is very fast)

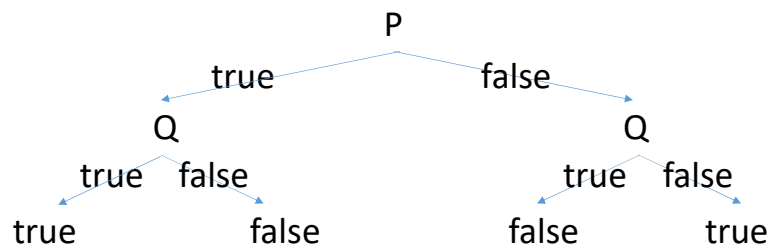
Feature_1	Feature_2	Class
false	false	? (guess: false)

## Example: Decision Tree (and Rule) Learning

- Can decision trees represent all Boolean functions?  
 $f(\text{Feature}_1, \dots, \text{Feature}_n) \equiv \text{some propositional sentence}$

## Example: Decision Tree (and Rule) Learning

- Can decision trees represent all Boolean functions? – Yes.  
 $f(\text{Feature}_1, \dots, \text{Feature}_n) \equiv \text{some propositional sentence}$
- Convert the propositional sentence into disjunctive normal form:  
 Example:  $(P \text{ AND } Q) \text{ OR } (\text{NOT } P \text{ AND } \text{NOT } Q)$



## Example: Decision Tree (and Rule) Learning

- Function learning needs bias, i.e. to prefer some functions over others.
- Occam's razor: "Small is beautiful."
- Here: Prefer small decision trees over large ones (e.g. with respect to their depth, their number of nodes, or (used here) **their average number of feature tests to determine the class**).
- Reason: ~~The functions encountered in the real world are often simple.~~
- Real reason: There are fewer small decision trees than large ones. Thus, there is only a small chance that ANY small decision tree that does not represent the correct function is consistent with all labeled examples.
- Problem: Finding the smallest decision tree that is consistent with all labeled examples is NP-hard. So, we just **try** to find a small decision tree.

## Example: Decision Tree (and Rule) Learning

- Real reason: There are fewer small decision trees than large ones. Thus, there is only a small chance that ANY small decision tree that does not represent the correct function is consistent with all labeled examples.
- In a country with 10 cities, if the majority of the population of a city voted for the winning president in the past 10 elections, perhaps they represent the "average citizen" of the country well.
- In a country with 10,000 cities, if the majority of the population of a city voted for the winning president in the past 10 elections, it could just be by chance. For example, if every citizen voted randomly for one of two candidates in the past 10 elections, there is still a good chance that there exists a city where the majority of the population voted for the winning president in the past 10 elections, just because there are so many cities.



## ID3 Algorithm

- The trivial decision trees (“always true” or “always false”) do not work here.

	Feature_1	Feature_2	Feature_3	Feature_4	Class
E(xample) 1	true	true	false	true	true
E(xample) 2	true	false	false	false	true
E(xample) 3	true	true	true	true	false
E(xample) 4	true	true	true	false	false

~~true~~

This decision tree does not work here since the examples do not all have class true.

~~false~~

This decision tree does not work here since the examples do not all have class false.

## ID3 Algorithm

- Put the most discriminating feature at the root.

	Feature_1	Feature_2	Feature_3	Feature_4	Class
E(xample) 1	true	true	false	true	true
E(xample) 2	true	false	false	false	true
E(xample) 3	true	true	true	true	false
E(xample) 4	true	true	true	false	false

~~Feature\_1  
true false  
E1: true  
E2: true  
E3: false  
E4: false~~

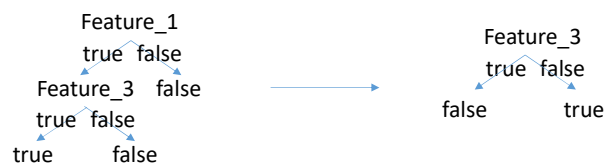
Feature\_2  
true false  
E1: true E2: true  
E3: false  
E4: false

Feature\_3  
true false  
E3: false E1: true  
E4: false E2: true

Feature\_4  
true false  
E1: true E2: true  
E3: false E4: false

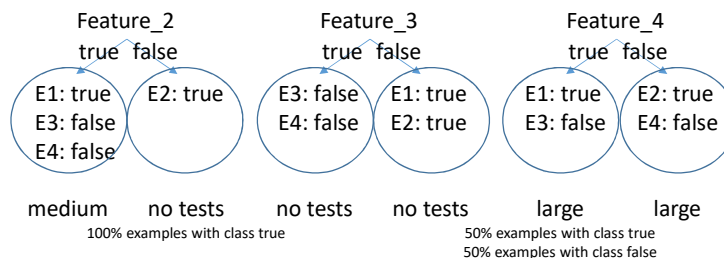
## ID3 Algorithm

- Putting Feature\_1 at the root is not helpful at all since all labeled examples have the same value for Feature\_1.
- If we eventually find a decision tree that is consistent with all labeled examples, then we can decrease the average number of feature tests to determine the class by deleting the root.



## ID3 Algorithm

- What's the average number of feature tests to determine the class in the following cases:
  - 1,00 (= 100%) examples with class true – no feature tests!
  - 999 examples with class true, 1 example with class false – likely: a very small number
  - 500 (= 50%) examples with class true, 500 (= 50%) examples with class false – likely: a large number
  - 1 example with class true, 999 examples with class false – likely: a very small number
  - 1,000 examples with class false – no feature tests!

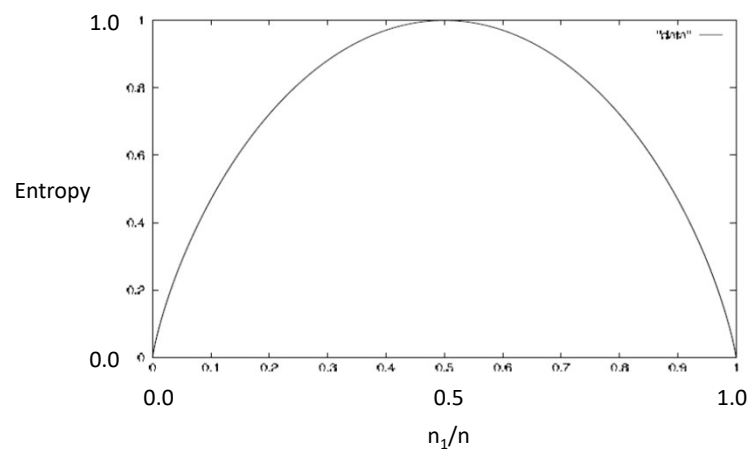


## ID3 Algorithm

- We use the entropy as a measure that we assume to be proportional to the average number of feature tests to determine the class, which we are **trying** to minimize (without guarantees).
- Assume that there are  $n$  examples and that  $n_i$  examples have class  $i$ .
- Entropy =  $-\sum_i (n_i/n) \log_2 (n_i/n)$
- The entropy is zero if no feature tests are necessary.
- Remember that  $\log_2 x = \ln x / \ln 2 = \log_{10} x / \log_{10} 2$ .
- Remember that  $\lim_{x \rightarrow 0} (x \log_2 x) = 0$ . Thus, “ $0 \log_2 0 = 0$ .”

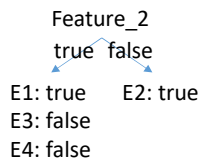
## ID3 Algorithm

- In our case, there are 2 classes.



## ID3 Algorithm

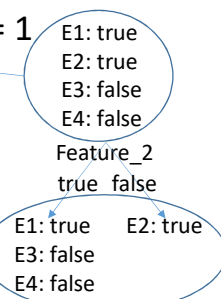
- Put the feature at the root that results in the smallest average entropy after splitting the examples.



- Left branch:
  - 3 out of 4 examples go down the left branch.
  - The entropy of the 3 examples is  $-(1/3 \log_2 (1/3) + 2/3 \log_2 (2/3)) = 0.9182$ .
- Right branch:
  - 1 out of 4 examples go down the right branch.
  - The entropy of the 1 example is  $-(1/1 \log_2 (1/1) + 0/1 \log_2 (0/1)) = 0$ .
- The average entropy after splitting the examples is  $\frac{3}{4} 0.9182 + \frac{1}{4} 0 = 0.6887$ .

## ID3 Algorithm

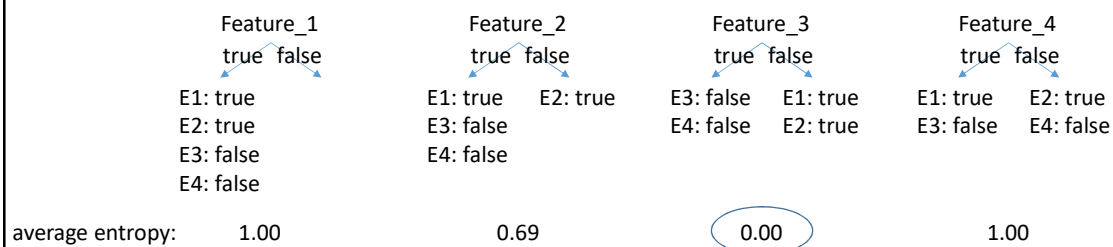
- The textbook does not pick the feature that results in the smallest average entropy. Instead, it (equivalently) picks the feature that results in the largest information gain, which is the entropy of the examples (= before splitting them) minus the average entropy after splitting them.
- The entropy of the examples is  $-(2/4 \log_2 2/4 + 2/4 \log_2 2/4) = 1$
- The average entropy after splitting them is 0.6887.
- The information gain is  $1 - 0.6887 = 0.3113$



## ID3 Algorithm

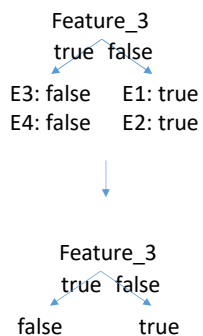
- Put the feature at the root that results in the smallest average entropy.

	Feature_1	Feature_2	Feature_3	Feature_4	Class
E(xample) 1	true	true	false	true	true
E(xample) 2	true	false	false	false	true
E(xample) 3	true	true	true	true	false
E(xample) 4	true	true	true	false	false



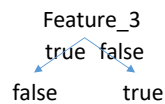
## ID3 Algorithm

- Put the feature at the root that results in the smallest average entropy.



## ID3 Algorithm: Complete Example

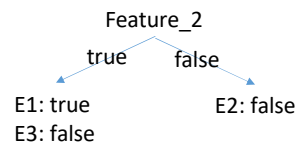
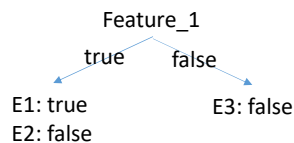
	Feature_1	Feature_2	Feature_3	Feature_4	Class
E(xample) 1	true	true	false	true	true
E(xample) 2	true	false	false	false	true
E(xample) 3	true	true	true	true	false
E(xample) 4	true	true	true	false	false



Feature_1	Feature_2	Feature_3	Feature_4	Class
false	false	false	false	? (guess: true)
true	true	true	true	? (guess: false)

## ID3 Algorithm: Complete Example

	Feature_1	Feature_2	Class
E(xample) 1	true	true	true
E(xample) 2	true	false	false
E(xample) 3	false	true	false



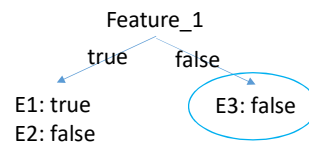
average entropy:  $2/3 \cdot 1 + 1/3 \cdot 0 = 2/3$

$2/3 \cdot 1 + 1/3 \cdot 0 = 2/3$

We have a tie that we can break arbitrarily.

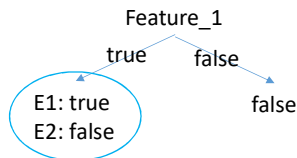
## ID3 Algorithm: Complete Example

	Feature_1	Feature_2	Class
E(xample) 1	true	true	true
E(xample) 2	true	false	false
E(xample) 3	false	true	false



Now, we solve this classification problem recursively.

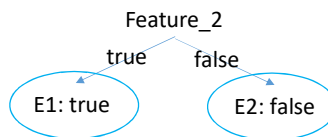
Now, we solve this classification problem recursively.



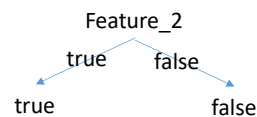
## ID3 Algorithm: Complete Example

	Feature_1	Feature_2	Class
E(xample) 1	true	true	true
E(xample) 2	true	false	false

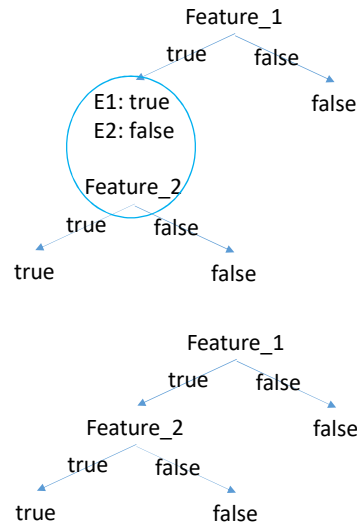
Now, we solve this classification problem recursively.



Now, we solve this classification problem recursively.

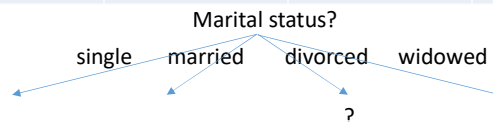


## ID3 Algorithm: Complete Example



## Issues with Multi-Valued Features, Issues with Missing Feature Values

How old are they?	What is their current salary per year?	Do they have a savings account?	What is their marital status?	...	Would you issue a credit card to them?
52	\$150,000	yes	widowed	...	yes
40	\$50,000	no	married	...	no
20	\$60,000	yes	married	...	yes
31	\$20,000	yes	single	...	yes



How old are they?	What is their current salary per year?	Do they have a savings account?	What is their marital status?	...	Would you issue a credit card to them?
26	\$40,000	no	divorced	...	?



## Example: Decision Tree (and Rule) Learning

- Want to play around with decision tree learning?
- Go here: <http://aispace.org/dTree/>