

# Supervised learning - Classification

## Decision Trees

9 Yes  
6 No

ID	Age	Has_job	Own_house	Credit_rating	Class
1	young	false	false	fair	No
2	young	false	false	good	No
3	young	true	false	good	Yes
4	young	true	true	fair	Yes
5	young	false	false	fair	No
6	middle	false	false	fair	No
7	middle	false	false	good	No
8	middle	true	true	good	Yes
9	middle	false	true	excellent	Yes
10	middle	false	true	excellent	Yes
11	old	false	true	excellent	Yes
12	old	false	true	good	Yes
13	old	true	false	good	Yes
14	old	true	false	excellent	Yes
15	old	false	false	fair	No

Many definitions of "smallest"

Bad News For any reasonable definition, NP complete

Heuristic ("greedy")

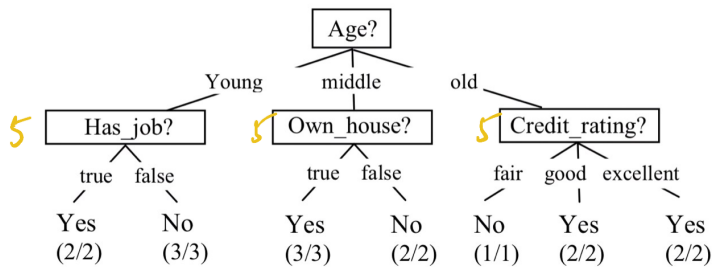
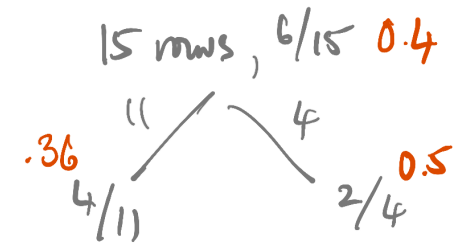
Which attribute to query next?

Impurity of data set

$$\frac{\text{Minority \#}}{\text{Total}}$$

Reduce impurity

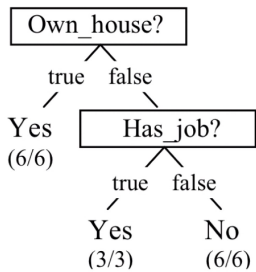
Weighted avg



Prefer small trees

- Occam's Razor

- Minimize overfitting



Better ways to measure impurity

Information theory [Shannon]

Messages over 4 letter alphabet {a,b,c,d}

Encode in binary

Natural choice is 2 bits/letter

a 00  
b 01  
c 10  
d 11

101100101100

N chars → 2N bits

26 letters a-z 5 bits/letter

Variable length coding

More frequent - shorter code

Characters  $c_1, c_2, \dots, c_k$

Frequency/prob  $p_1, p_2, \dots, p_k$

$$\sum p_i = 1$$

$$\sum_i -p_i \log_2 p_i = \text{Shannon entropy}$$

$$\sum -p_i \log_2 p_i$$

{a, b, c, d}

$\frac{1}{4} \frac{1}{4} \frac{1}{4} \frac{1}{4}$

$$\sum_{i=1}^4 -\frac{1}{4} \log_2 \frac{1}{4}$$

$$\log_2 \frac{1}{4}$$

$$= \log_2 (4)^{-1}$$

$$= -1 \cdot \log_2 4$$

$$= -2$$

$$\sum_{i=1}^4 -\frac{1}{4} \cdot 2$$

$$= \sum_{i=1}^4 \frac{1}{2} = 2$$

a b c d  
 $\frac{1}{2} \frac{1}{4} \frac{1}{8} \frac{1}{8}$

$$\frac{1}{2} \log \frac{1}{2} + \frac{1}{4} \log \frac{1}{4} + \frac{1}{4} \log \frac{1}{8}$$

$$\frac{1}{2} + \frac{1}{2} + \frac{3}{4}$$

$$= 1 \frac{3}{4}$$

Entropy is maximized when each  $p_i = \frac{1}{k}$

k=2

$$p_1 = p_2 = \frac{1}{2}$$

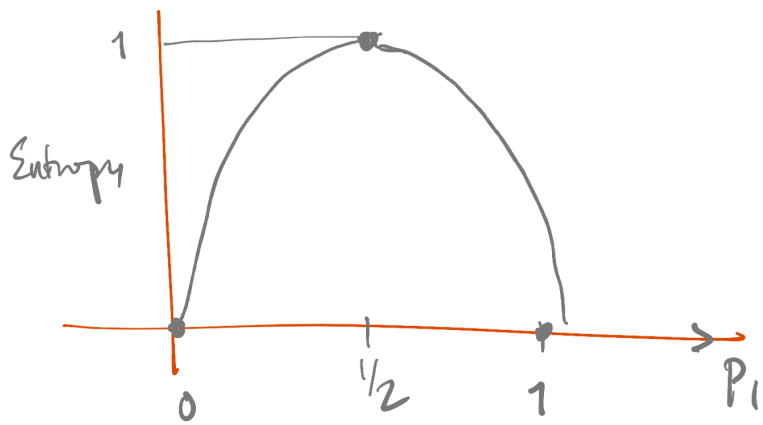
$$\sum_{i=1}^2 -\frac{1}{2} \log \frac{1}{2}$$

$$= 1$$

$$p_1 = 0 \quad p_2 = 1$$

$$\frac{0 \log 0}{= 0}$$

$$\frac{1 \log 1}{= 0}$$



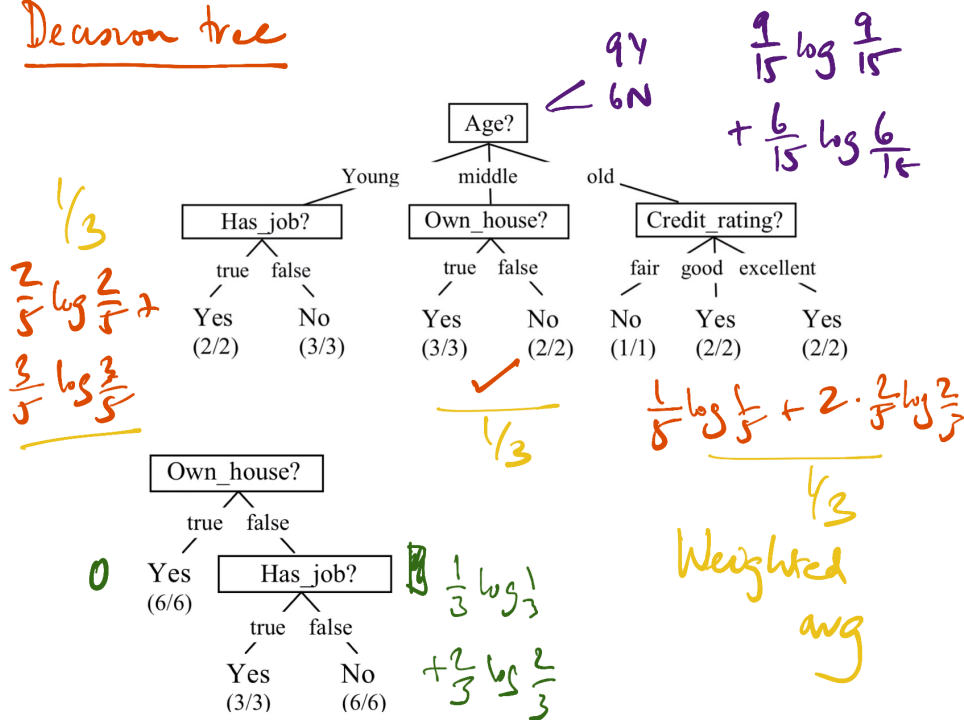
$$\sum_{i=1}^k \frac{1}{k} \log \frac{1}{k} \quad \log_2 k, k > 2$$

Choose attribute that reduces entropy the most **Information Gain**

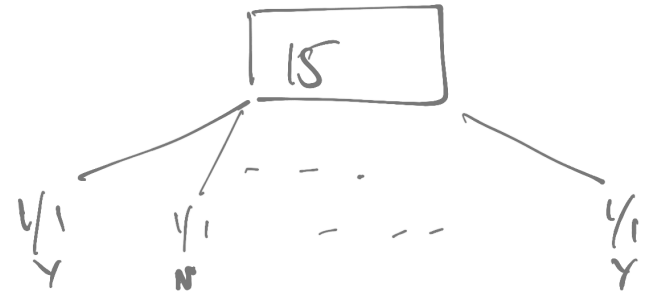
## Ross Quinlan

Now, this description does violence to statistical notions of sampling and confidence limits, so the reasoning should be taken with a large grain of salt. Like many heuristics with questionable underpinnings, however, the estimates that it produces seem frequently to yield acceptable results.

## Decision tree



What if Aadhaan is an attribute?



$$\text{Entropy of Aadhaan} = \sum_{i=1}^k \frac{1}{k} \log \frac{1}{k}$$

Normalize information gain by entropy of attribute

Information Gain = Reduction in Entropy

Entropy of Attribute

Algorithm

Build Tree ( $A_1, \dots, A_k$ )

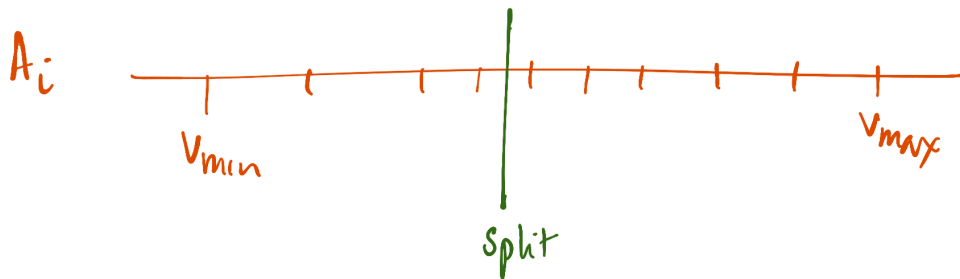
Choose  $A_j$  that maximizes Info Gain Ratio

For each child created

Build Tree ( $A_1, \dots, A_{j-1}, A_{j+1}, \dots, A_k$ )

Stop if  $\left\{ \begin{array}{l} \text{homogeneous} \\ \text{no attributes left} \end{array} \right. \rightarrow$  use majority

What about numeric attributes



$A_i < \text{split} ?$

Testing a classifier

Set aside some training data for validation

Randomly split as  $\left\{ \begin{array}{l} \text{Training} \\ \text{Test} \end{array} \right.$

Cross validation

10 x remove 10%

What is correctness?

Percentage of correct answers

Accuracy

Most classifiers are looking at skewed samples

Vast majority  $\rightarrow$  "Normal"

Small minority  $\rightarrow$  "Abnormal"

Refined calculation

Classifies says

	Y	N
Actual data	Y	N
	0	5
	0	95

Yes is the Abnormal (minority) case

Usually Precision ↑ Recall ↓

Selection of candidates

Screening Exam

High Recall

+

Interview

High Precision

Precision

	25	25
Recall	50	900

$$\frac{25}{50}$$

$$\frac{50}{75} \text{ mistakes}$$

$$\frac{25}{75}$$