

DMML, 23 Jan 2019

## Decision Trees

Items: Attributes  $A_1, A_2, \dots, A_n$ ,  
Category  $C$

Training data  $\rightarrow$  Model

Given  $(a_1, a_2, \dots, a_n)$ , predict  $c \in C$

Building a model - uniform algorithm

Specific training data  $\mapsto$  Specific model

## Decision trees

- Ask questions about attribute values
  - Adaptive - next question depends on previous answers
- Each question prunes available possibilities
- Stop when
  - Current set has a uniform class
  - Exhausted set of attributes to query
    - Answer is majority value

# Which question to ask?

- Prefer small trees
- Exactly computing smallest tree is NP complete
- Instead, greedy heuristic

Measure: "Impurity"  $\sim$  uncertainty

Naive measure: Error Rate      Minority %

More sophisticated measures of impurity/uncertainty

Information theory

$\{a, b, c, d, \dots, z\}$

Send a string of  
chars in binary

Uniform size encoding = 5 bits/char

Message of  $N$  chars  $\rightarrow 5N$  bits

In fact - letters do not occur uniformly

Use variable length encoding - shorter seq  
for freq letters

# Claude Shannon

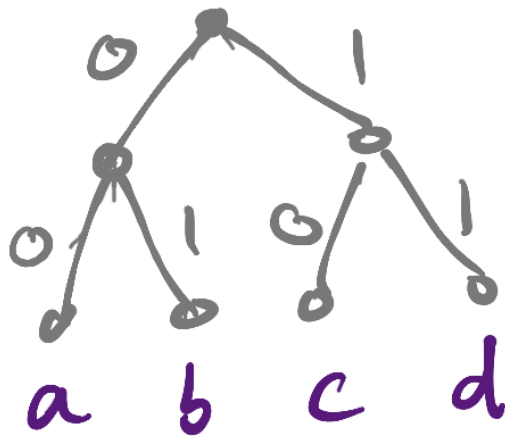
$\{a, b, c, d\}$

$\begin{array}{cccc} | & | & | & | \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{8} & \frac{1}{8} \end{array}$

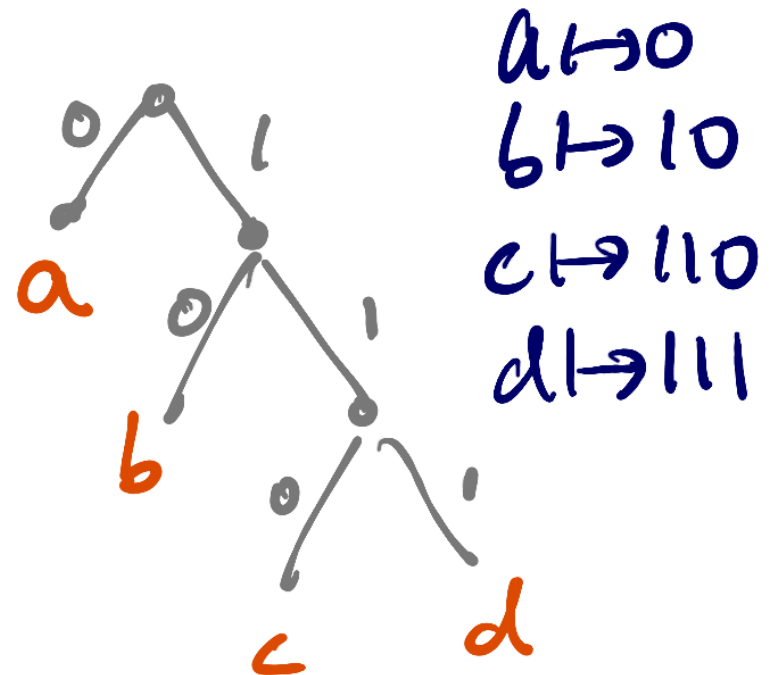
— uniform 2 bit code

$N \text{ char} \rightarrow 2N \text{ bits}$

Uniform encoding



Non uniform



200 characters

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

Uniform = 400 bits

Nonuniform =

100 a's	50 b's	25 c's	25 d's
↓	↓	↓	↓
100	+ 100	+ 75	+ 75
= 350			

Reverse analysis  $\Rightarrow$  1.75 bits/char

Is this optimal?

Distribution of characters reflects  
uncertainty in their values

## Shannon entropy

Letters/outcomes

$a_1$   $a_2$  ...  $a_k$

Probabilities

$p_1$   $p_2$  ...  $p_k$

$$\sum p_i = 1$$

$$\text{Entropy} \stackrel{\text{defn}}{=} - \sum_{i=1}^k p_i \log_2 p_i$$

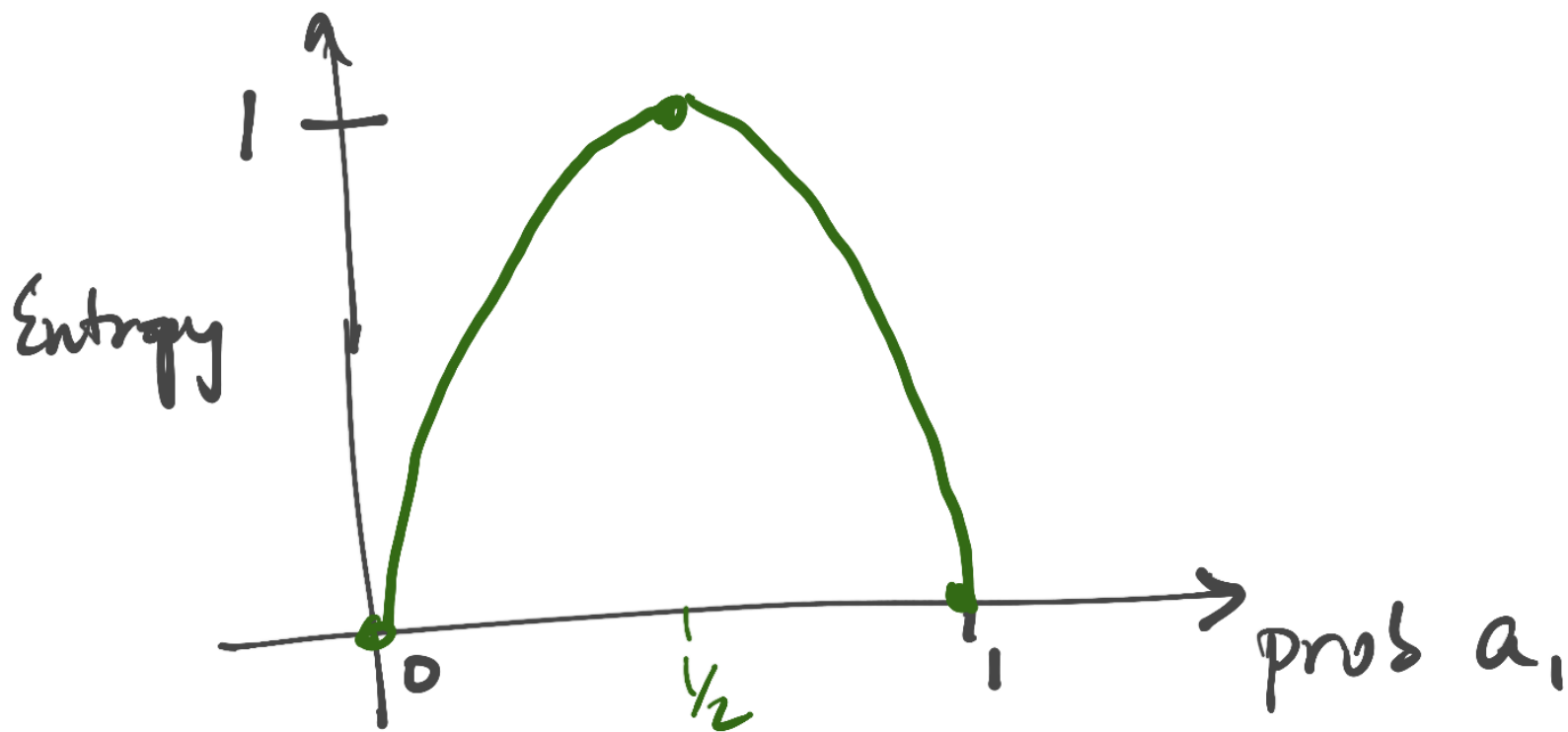
Entropy measures uncertainty

$$\begin{matrix} a_1 & a_2 \\ 1/2 & 1/2 \end{matrix}$$

$$- \left( \underbrace{\frac{1}{2} \log \frac{1}{2}}_{-\frac{1}{2}} + \underbrace{\frac{1}{2} \log \frac{1}{2}}_{-\frac{1}{2}} \right) = 1$$

$$\begin{matrix} 0 & 1 \\ 1 & 0 \end{matrix} \left. \vphantom{\begin{matrix} 0 & 1 \\ 1 & 0 \end{matrix}} \right\} \rightarrow = 0 \quad \left( \underbrace{0 \log 0}_{\text{assume } 0} + \underbrace{1 \log 1}_0 \right)$$



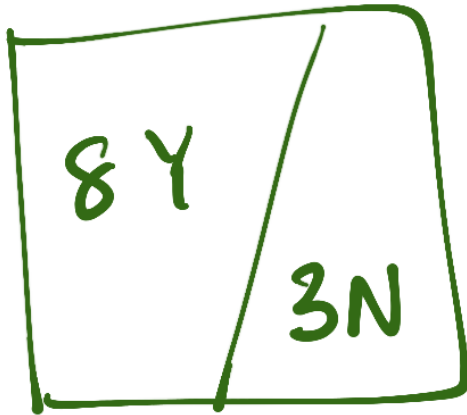


For  $a_1, \dots, a_k$  max entropy is  $P_i = \frac{1}{k}$

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

$$\begin{aligned}
 & \frac{1}{2} \log \frac{1}{2} + \frac{1}{4} \log \frac{1}{4} \\
 & \quad + 2 \cdot \frac{1}{8} \log \frac{1}{8} \\
 & - (-1.75) = 1.75
 \end{aligned}$$

# Entropy for impurity



$$\text{Error Rate} = \frac{3}{11}$$

$$\text{Entropy} = - \left( \frac{8}{11} \log \frac{8}{11} + \frac{3}{11} \log \frac{3}{11} \right)$$

Ross Quinlan — Entropy works better than other notions on many benchmarks

C4.5

# Gini Index (Economics)

$$1 - \sum p_i^2$$

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

$$G.I = 1 - \frac{1}{2} = \frac{1}{2}$$

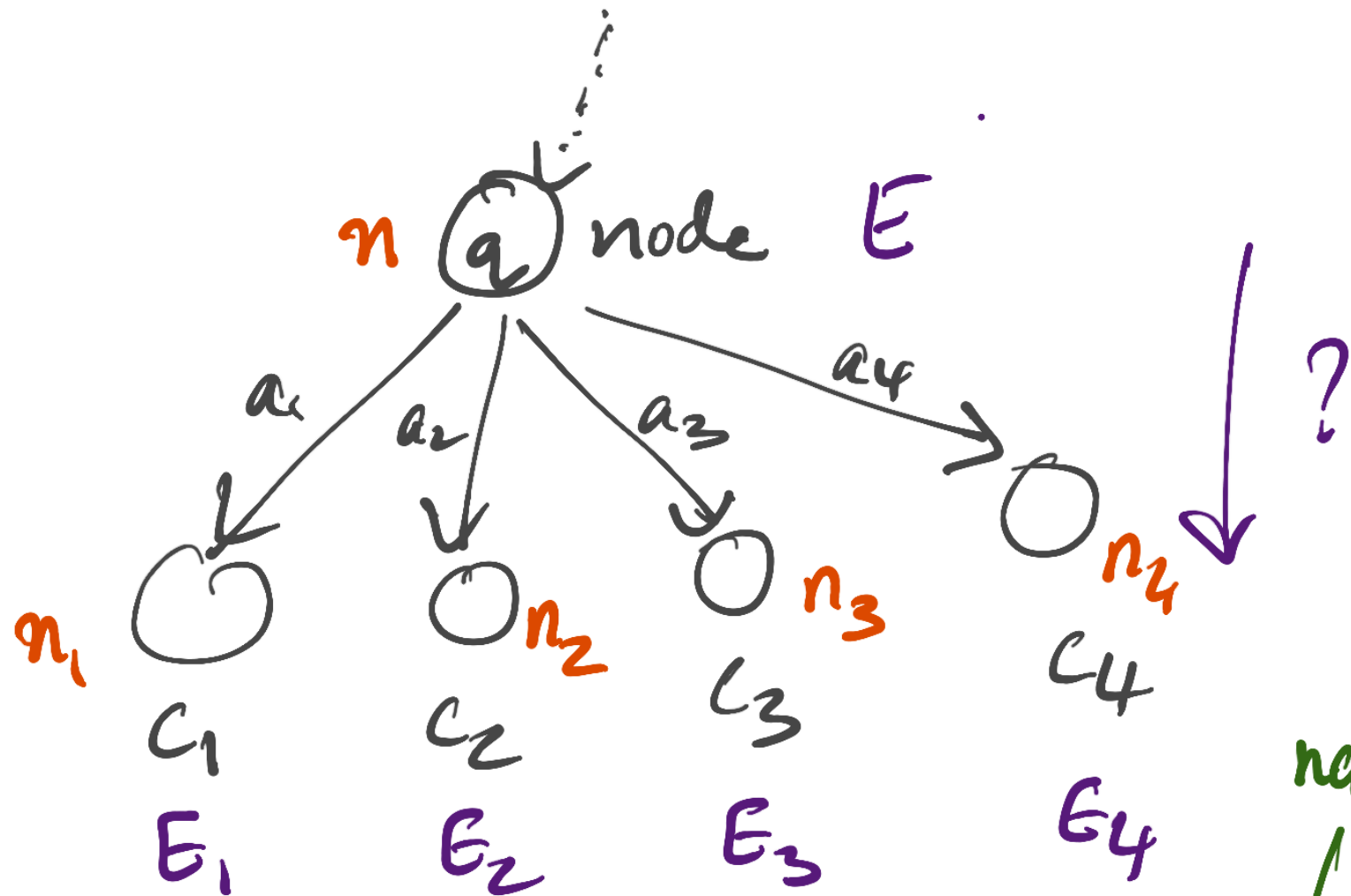
$$p_1 = p_2 = \dots = p_k = \frac{1}{k}$$

$$1 - k \cdot \frac{1}{k^2}$$

Classification And Regression Tree  
CART implemented in R

$$= \frac{k-1}{k}$$

Reduction in Impurity = Information Gain



$$n_1 + n_2 + n_3 + n_4 = n$$

Weighted  
Avg

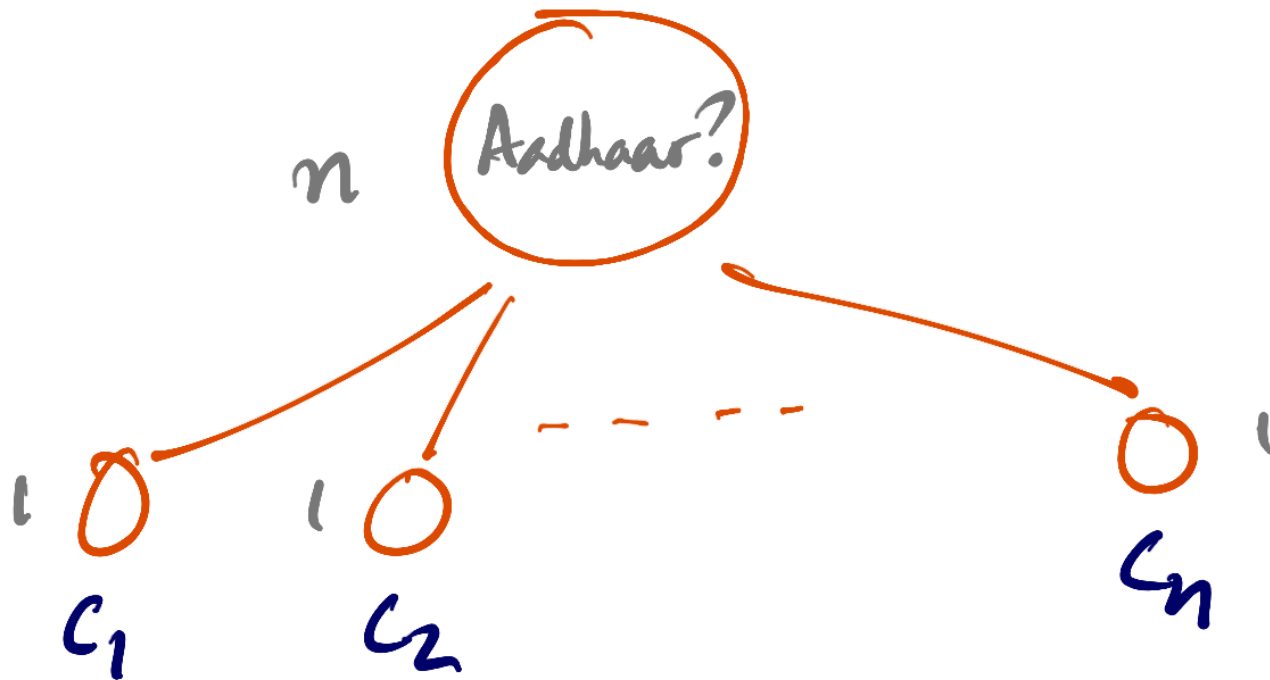
new  $E'$   
after  
 $q$

$$\sum_i \frac{n_i}{n} \cdot E_i$$

Special Case

What if  $A_i = \text{Aadhaar No?}$

Split on Aadhaar



$E_i = 0$  for all  $c_i \rightarrow \text{wt avg} = 0$

Measure entropy of an attribute

$$\text{Entropy} = k \log k - \sum \frac{1}{k} \log \frac{1}{k}$$

Moderate Information Gain (Absolute)  
by entropy of attribute

$$\text{Information Gain Ratio} = \frac{\text{Absolute Entropy Gain}}{\text{Entropy of Attribute}}$$

Borrowed Entropy from Info Theory

Why does it make sense?

Because it works!