

DMML, 16 Jan 2020

Decision Trees

Items with attributes

(a_1, \dots, a_n) & class c

Training data \rightarrow Model

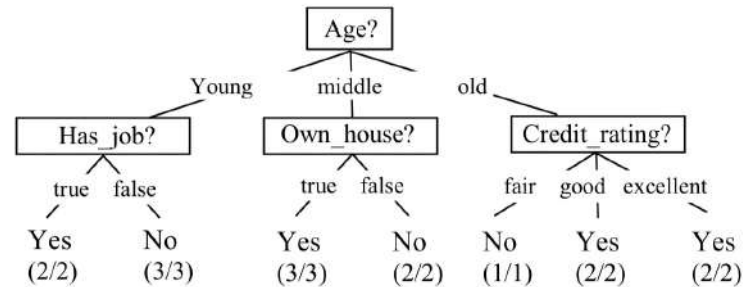
Adaptively ask question

Prefer small trees

Computationally hard

Greedy heuristic

| 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 |



Greedy heuristic

Reduce "impurity" as much as possible

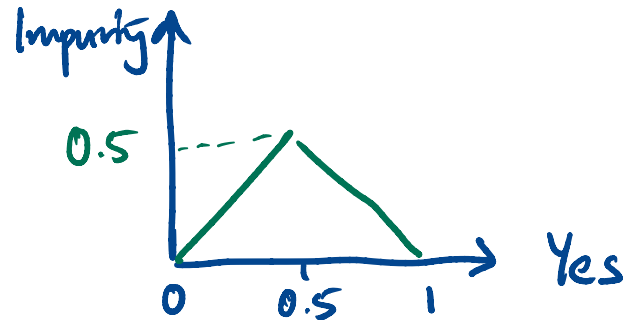
"Pure" = all items in the set have same class

Impure set - verdict is majority

Minority case = error = misclassification

Impurity = Misclassification Rate

Max impurity = 0.5

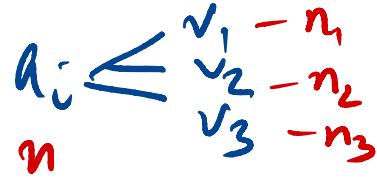
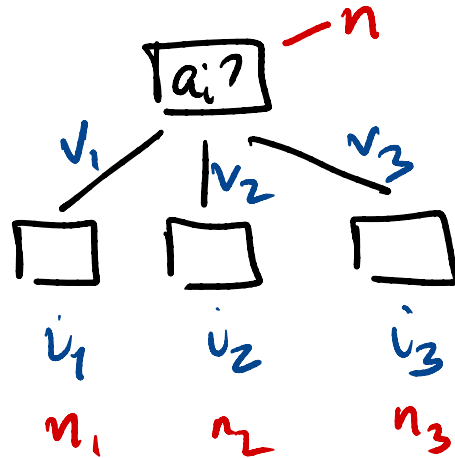


Tree building algorithm

Available attributes (a_1, \dots, a_k)

Evaluate current misclassification rate

For each a_i



Impurity

Wted
avg

$$\frac{n_1}{n} \cdot i_1 + \frac{n_2}{n} \cdot i_2 + \frac{n_3}{n} \cdot i_3$$

Compute weighted avg of impurity for each a_i

Among these, choose the best one $\rightarrow a_j$

For each child node, apply the same algorithm,

with $(a_1, \dots, a_{j-1}, a_{j+1}, \dots, a_n)$ as available attributes

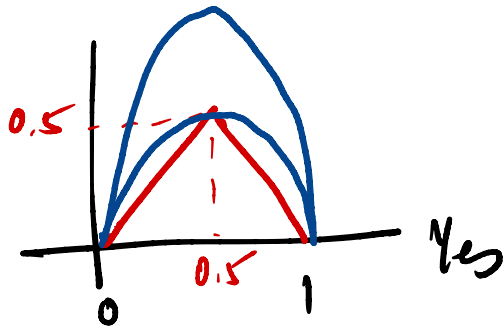
Stopping?

1. No more available attributes
2. Node is pure

1. Impurity = Error = "Cost"

Incrementally reducing cost by refining the model

2. Is there a better notion of impurity?



Better to have a measure that penalizes impurity non-linearly

Information Theory [Shannon]

Message alphabet \rightarrow Encoding \rightarrow Transmission Alphabet

$\{a, b, \dots, y, z\}$

26

$\{0, 1\}$

$2^4 < 26 \leq 2^5$

Can we use frequency information to improve this?

Interpret frequencies as probabilities

Example

{a, b, c, d}

$\xrightarrow{\text{encode}}$
{0,1}

2 bits

Message length N

$\longrightarrow 2N$

a - $\frac{1}{2}$

b - $\frac{1}{4}$

c - $\frac{1}{8}$

d - $\frac{1}{8}$

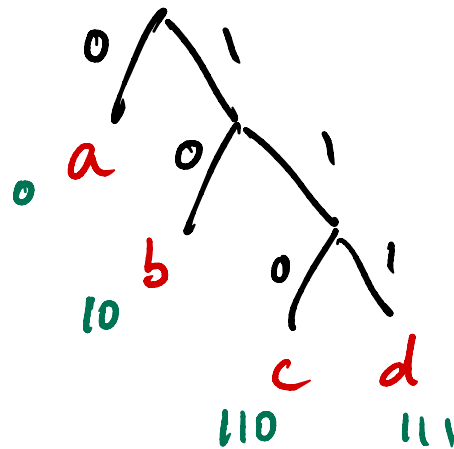
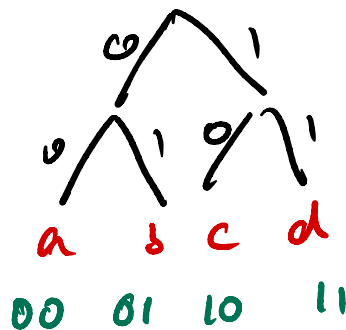
Variable length

encoding -

shorter codes

for frequent

letters



200 letters input $\xrightarrow{\text{Naive}}$ 400 bits

$$\begin{aligned} & | \\ & 100 a \times 1 + 50 b \times 2 + 25 c \times 3 + 25 d \times 3 \\ & \quad 100 + 100 + 75 + 75 = \underline{\underline{350}} \end{aligned}$$

Huffman Coding

Shannon - entropy $-\sum p_i \log_2 p_i$

p_i = probability of i th letter

$$\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8}$$

$$- \left(\frac{1}{2} \log \frac{1}{2} + \frac{1}{4} \log \frac{1}{4} + \frac{1}{8} \log \frac{1}{8} + \frac{1}{8} \log \frac{1}{8} \right)$$

$-\frac{1}{2} \quad + \quad -\frac{1}{2} \quad + \quad -\frac{3}{8} \quad + \quad -\frac{3}{8}$

$$- \left(-1 + \left(-\frac{3}{4}\right) \right)$$

$$- \left(-1.75 \right)$$

$= 1.75 \implies$ avg length per letter

is 1.75 bits (lower 3d)

200 length input $\times 1.75 = 350$ chars

Borrow entropy as a measure of impurity

$$N \text{ cases} \begin{cases} M \text{ Yes} \\ N-M \text{ No} \end{cases} \quad \begin{aligned} - P_Y &= \frac{M}{N} \\ - P_N &= \frac{N-M}{N} \end{aligned}$$

$$- (P_Y \log P_Y + P_N \log P_N)$$

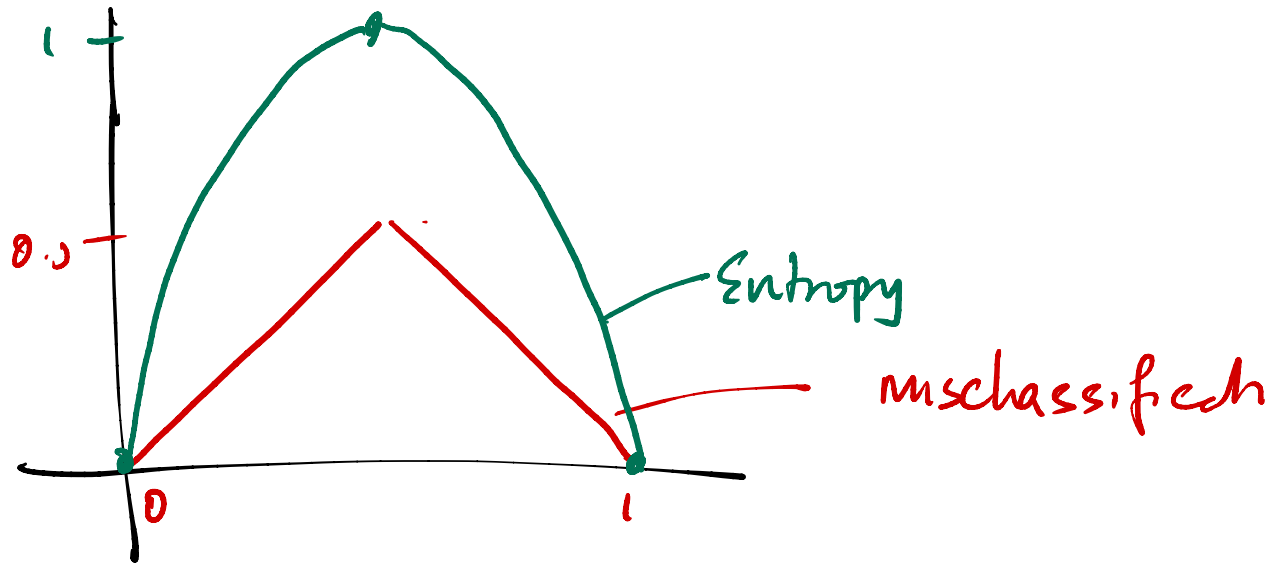
$$P_Y = 1, P_N = 0 \quad - \left(\underbrace{1 \log 1}_0 + 0 \underbrace{\log 0}_? \right) = 0$$

$$P_Y = 0, P_N = 1$$

For entropy

$0 \log 0$
defined = 0

$$P_Y = P_N = \frac{1}{2} \quad - \left(\frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right)$$
$$= - \left(-\frac{1}{2} + -\frac{1}{2} \right) = +1$$



Another option

Economics - Equality of distribution of wealth/resources

$$1 - \sum p_i^2$$

$$P_Y = 1 \quad P_N = 0$$

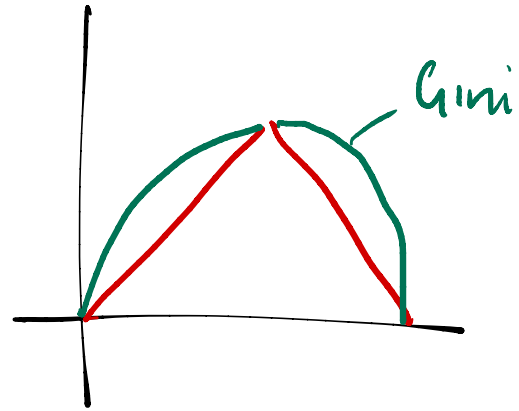
$$0$$

$$P_N = 1 \quad P_Y = 0$$

$$0$$

$$P_Y = P_N = \frac{1}{2}$$

$$1 - \left(\frac{1}{4} + \frac{1}{4}\right) = \frac{1}{2}$$



$$\frac{1}{3}, \frac{1}{3}, \frac{1}{3}$$

$$\underline{\text{Entropy}} = \left(\frac{1}{3} \log \frac{1}{3} \right) \cdot 3$$

$$\log_2 3$$

$$\underline{\text{Gini}} = 1 - \sum \left(\frac{1}{3} \right)^2 = \frac{2}{3}$$

Most implementations use Gini index

Any such function should do

Shannon

Information is inverse of entropy

For us entropy = impurity

Reduce impurity = Gain information