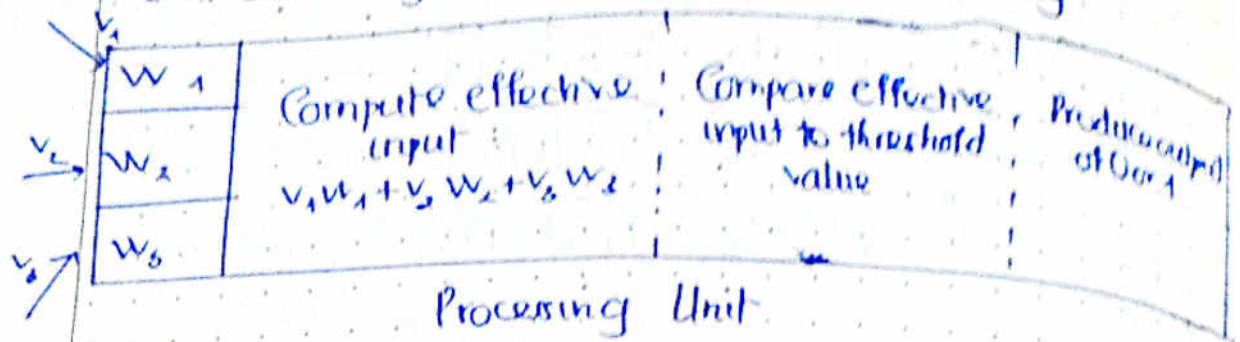


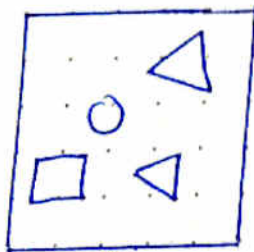
nhận dạng mẫu

B. Hoạt động bên trong của một đơn vị xử lý

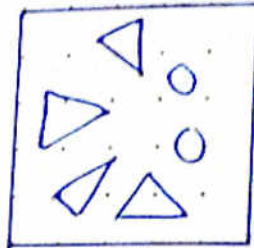


## HỌC QUA LOGIC

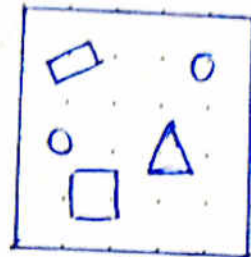
Lớp B



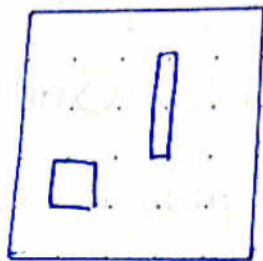
(6)



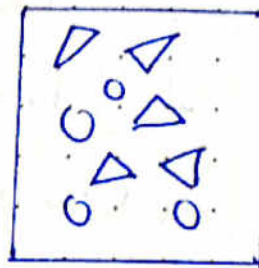
(7)



(8)



(9)



(10)

P1:  $\exists$  hình tam giácP2:  $\exists$  hình tròn

Lớp	Hình	Tam giác (P1)	Hình tròn (P2)	Tròn (P3)	Chữ nhật (P4)	Đa giác (P5)
B	6	1	1	0	1	0
B	7	1	1	0	0	0
B	8	1	1	0	1	0
B	9	0	0	0	1	0
B	10	1	0	1	0	0

$$q_1 = \bar{p}_1 p_2 \bar{p}_3 p_4 \bar{p}_5$$

$$q_2 = p_1 p_2 \bar{p}_3 \bar{p}_4 \bar{p}_5$$

$$q_3 = p_1 p_2 \bar{p}_3 p_4 \bar{p}_5$$

$$q_4 = \bar{p}_1 \bar{p}_2 \bar{p}_3 p_4 \bar{p}_5$$

$$q_5 = p_1 \bar{p}_2 p_3 \bar{p}_4 \bar{p}_5$$

$$p_1 p_2 \bar{p}_3 \bar{p}_4 \bar{p}_5 \vee \dots$$

$$p_1 \bar{p}_2 \bar{p}_3 p_4 \bar{p}_5 \vee p_1 p_2 \bar{p}_3 \bar{p}_4 \bar{p}_5$$

$$p_1 \bar{p}_2 p_3 p_4 \bar{p}_5$$

$$\bar{p}_3 p_3 \bar{p}_5$$

$$X \in D : X \text{ thỏa}$$

$$Q = p_1 p_2 \bar{p}_3 p_4 \bar{p}_5 \vee p_1 p_2 \bar{p}_3 \bar{p}_4 \bar{p}_5 \vee p_1 p_2 \bar{p}_3 p_4 \bar{p}_5 \\ \vee \bar{p}_1 \bar{p}_2 \bar{p}_3 p_4 \bar{p}_5 \vee p_1 \bar{p}_2 p_3 \bar{p}_4 \bar{p}_5$$

$$X \in D : X \in \left[ \begin{array}{l} \dots \\ \dots \end{array} \right]$$

## II. Học qua tư quan sát

### 1. Ý tưởng thuật toán Quinlan:

- + Cho một bảng quan sát là tập hợp các mẫu và các thuộc tính nhất định của các đối tượng nào đó.
- + Sử dụng một độ đo để định lượng và đề ra tiêu chuẩn nhằm chọn lựa một thuộc tính mang tính chất "phân loại" để phân bảng thành các hàng con nhỏ hơn, rồi cho từ mỗi bảng con này dễ dàng phân tích tìm ra một quy luật chung



(Định quy)

+ Tổ đề thiết lập dưới "Cây quyết định" cho thấy thứ tự của các thuộc tính đang xét.

VD: Xác định là  $\bar{A}$  châu Á hay châu Âu khi xem xét một nhóm người cần đi học: hình dáng, chiều cao giới tính.

- Định nghĩa đồ cho V:

$$+ V(H, \text{Dáng} = T_0) = (A'_{T_0}, A''_{T_0})$$

$$A'_{T_0} = \frac{\text{Tổng số quan sát châu Á có } (H, \text{Dáng} = T_0)}{\text{Tổng số quan sát có } (H, \text{Dáng} = T_0)}$$

$$A''_{T_0} = \frac{\text{Tổng số quan sát châu Âu có } (H, \text{Dáng} = T_0)}{\text{Tổng số quan sát có } (H, \text{Dáng} = T_0)}$$

- Tiêu chuẩn phân loại: Chọn thuộc tính có nhiều vector rời rạc nhất.

Giả sử có bảng quan sát sau:

STT	Hình dáng	Chiều cao	Giới tính	Quan sát
1	Tô	Trung bình	Nam	Châu Á
2	Nhỏ	Thấp	Nam	Châu Á
3	Nhỏ	Trung bình	Nam	Châu Á
4	Tô	Cao	Nam	Châu Âu
5	Nhỏ	Trung bình	Nữ	Châu Âu
6	Nhỏ	Cao	Nam	Châu Âu
7	Nhỏ	Cao	Nữ	Châu Âu
8	Tô	Trung bình	Nữ	Châu Âu

- Tình trạng:

+ Lần 1:

• Hình dáng:

$$V(H. \text{Dạng} = \text{to}) = (1/3, 2/3)$$

$$V(H. \text{Dạng} = \text{nhỏ}) = (2/5, 3/5)$$

• Chiều cao:

$$V(\text{Chiều cao} = \text{TB}) = (1/2, 1/2)$$

$$V(\text{Chiều cao} = \text{Thấp}) = (1, 0)^*$$

$$V(\text{Chiều cao} = \text{Cao}) = (0, 1)^*$$

• Giới tính:

$$V(\text{Giới tính} = \text{Nam}) = (3/5, 2/5)$$

$$V(\text{Giới tính} = \text{Nữ}) = (0, 1)$$

STT	Hình dáng	Giới tính	Quan sát
1	To	Nam	Châu A
3	Nhỏ	Nam	Châu A'
5	Nhỏ	Nữ	Châu A''
8	To	Nữ	Châu A'''

+ Lần 2:

• Hình dáng:

$$V(H. \text{Dạng} = \text{To}) = (1/2, 1/2)$$

$$V(H. \text{Dạng} = \text{Nhỏ}) = (1/2, 1/2)$$



Giải hình:

$$V(G.Tính = Nam) = (1, 0)^*$$

$$V(G.Tính = Nữ) = (0, 1)^*$$

Luật:

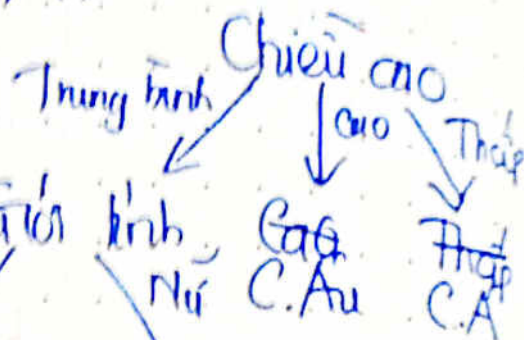
+ Nếu chiều cao = Cao

thì châu Âu

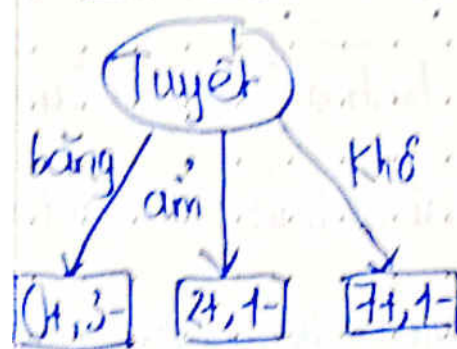
+ Nếu chiều cao = Thấp thì châu Á

+ Nếu chiều cao = T.bình và G.Tính = Nam thì châu Á

+ Nếu chiều cao = T.bình và G.Tính = Nữ thì châu Âu



2. Chọn thuộc tính lần 1

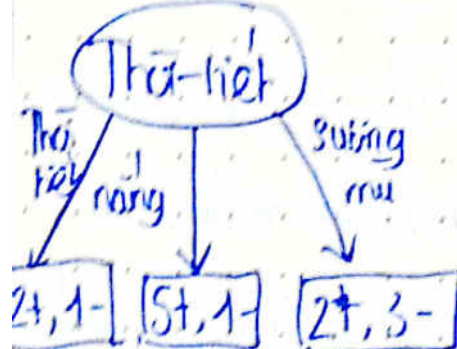


$$H_{băng} = -0/3 \log_2 0/3 - 3/3 \log_2 3/3 = 0$$

$$H_{ẩm} = -2/3 \log_2 2/3 - 1/3 \log_2 1/3 = 0,918$$

$$H_{khô} = -7/8 \log_2 7/8 - 1/8 \log_2 1/8 = 0,544$$

$$AE = 3/14 * 0 + 3/14 * 0,918 + 8/14 * 0,544 = 0,508$$



$$H_{thư} = -2/3 \log_2 2/3 - 1/3 \log_2 1/3 = 0,918$$

$$H_{rỗng} = -5/6 \log_2 5/6 - 1/6 \log_2 1/6 = 0,65$$

$$H_{rườ} = -2/5 \log_2 2/5 - 3/5 \log_2 3/5 = 0,971$$

$$AE = 3/14 * 0,918 + 6/14 * 0,65 + 0,971 * 5/14 = 0,822$$

### 3. Học là gì?

- Học là quá trình tiếp nhận và thức mới hoặc cấp nhất hi thức về hành vi, kĩ năng, ... và biến quan hệ thông tin  $\neq$  nhau.

+ Thay đổi để tốt hơn (Theo 1 điều kiện định trước)  
khi có hình huống tương tự xảy ra.

+ Học không phải là " học thuộc lòng "

- Học là một trong những khả năng quan trọng của con người.

### 4. Các mức độ của học

- Học là 1 quá trình được chia làm nhiều mức độ:

1. Được truyền tải hi thức và ghi nhớ chúng.

2. Tiếp nhận hi thức thông qua quan sát và tham dự các sự kiện.

3. Cải thiện hi thức thông qua luyện tập các kỹ năng vận động và nhận thức.

4. Tổ chức hi thức mới thành các biểu diễn tổng quát, hiệu quả.



## 5. Học máy là gì?

Học máy (Machine Learning) nghiên cứu các thay đổi mang tính thực nghiệm trong hệ thống, cho phép hệ thống thực thi những tác vụ tương tự các biểu thức logic qua nhiều lần. (Herbert Simon)

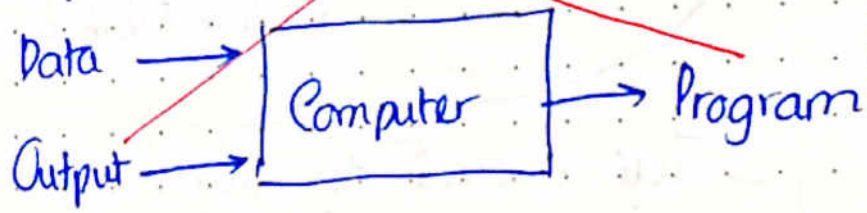
Learning - Pattern Detection - Data - Self - Programming

Nhiệm vụ của học máy là thiết kế các chương trình máy tính có thể học luật từ dữ liệu, thích nghi với thay đổi và cải thiện hiệu quả thông qua kinh nghiệm.

lập trình hệ thống



Học máy



Algorithms = Hạt giống

Data = Chất dinh dưỡng

You = Kiến trúc sư

Programs = Cây

## 6. Tại sao học lại khó?

Cho trước một lượng hữu hạn dữ liệu huấn luyện. Ta cần suy ra một quan hệ trong không gian vô hạn.

Trong thực tế, có vô số quan hệ như thế.

## 7. Nguyên lý Occam Razor

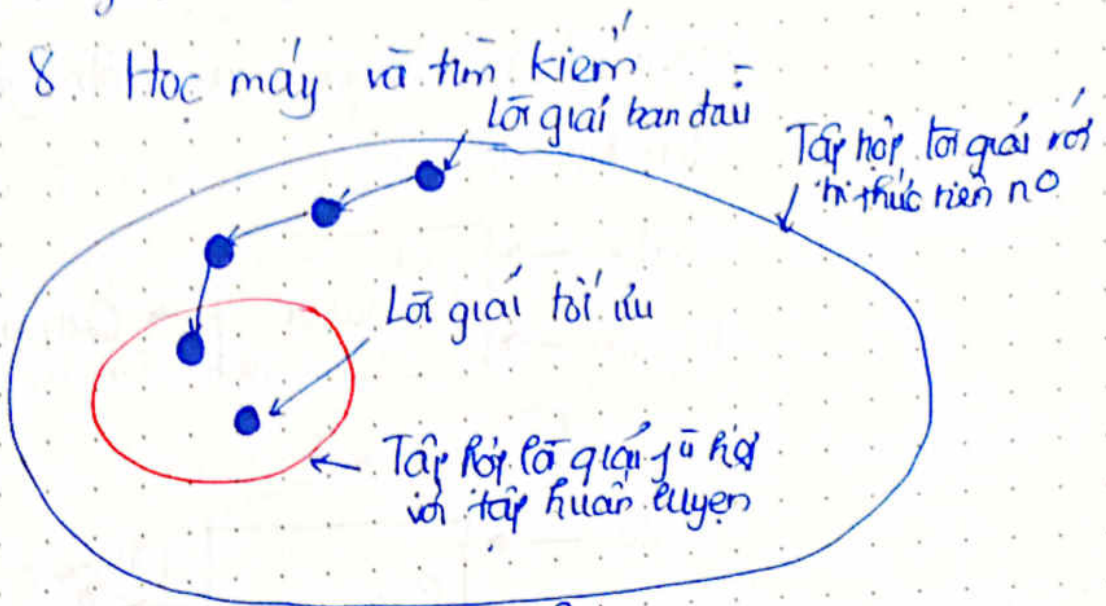
- Nguyên lý: Giải pháp đơn giản là giải pháp tốt.  
(William of Occam)

- Khi có nhiều giải pháp cho cùng một vấn đề, ta nên chọn giải pháp đơn giản nhất.

- Ta quan niệm thế nào là đơn giản?

+ Số ít thuộc tính  $n^0$  của bài toán để định nghĩa giải pháp đơn giản. Vd: Smoothness prior.

## 8. Học máy và tìm kiếm



- Tại sao cần học máy?

+ Sự phát triển của học máy nhờ vào 1 số ý tưởng sau:

- kiến thức mới về lý thuyết và thuật toán.

- Dữ liệu mức huấn luyện tăng trưởng như vũ bão.

- Sức mạnh tính toán đã sẵn sàng.



Ngành Công nghiệp phát triển nở rộ

9. Ứng dụng của học máy

- Ba lĩnh vực thích hợp cho học máy:

+ Khai thác dữ liệu: Sử dụng dữ liệu thu thập trong quá khứ để cải thiện quyết định.

+ Các phần mềm ứng dụng to thế lập trình bằng tay

+ Các chương trình tự điều chỉnh.

+ ...

- Xử lý thị giác:

+ Phát hiện / xác minh khuôn mặt

+ Nhận dạng chữ viết tay

- Xử lý tiếng nói: Nhận dạng âm vị / từ / câu / đối tượng.

- Các ứng dụng khác:

+ Đánh chỉ mục: Google, khai thác văn bản, truy vấn thông tin.

+ Tài chính: Dự đoán khả năng vay nợ, quản lý rủi ro đầu tư.

+ Truyền thông: Dự đoán tình hình giao thông.

+ Trò chơi: Cờ vua, cờ vây.

- + Điều khiển: Học máy Reinforcement
- Xã cần nhiều ứng dụng
  - Học máy được ứng dụng trong hướng hợp:
    - + Con h. có thể tham gia vào quá trình
    - + Con h. có thể giải thích về chuyên môn
    - + Các giải pháp thay đổi theo thời gian
    - + Các pháp cần giải thích giải thích nghi vs - tương tác cụ thể
  - Các lĩnh vực liên quan: Electrical engineering, Statistics, Philosophy, Neuroscience, Psychology, CS.
  - + Xác suất thống kê: Làm khớp mô hình với dữ liệu và kiểm tra kết quả
  - + Khai thác dữ liệu / Phân tích dữ liệu thăm dò: Phát hiện hình mẫu trong dữ liệu
  - + Lý thuyết điều khiển thích nghi: Học trực tuyến các mô hình và dùng chúng để đạt tới một số mục tiêu
  - + Trí tuệ nhân tạo: Xây dựng máy thông minh

## 10. Các dạng học máy

### 10.1. Học có giám sát (Supervised learning)



- Dựa vào các mẫu huấn luyện đã biết phân lớp
- Bộ học chỉ hướng dẫn để nhận diện các mẫu cho mỗi phân lớp.

## 10.2 Học không giám sát (Unsupervised learning)

- Dựa vào mẫu huấn luyện chưa biết phân lớp.
- Bộ học cần phải hệ hàm kiểm tra tổng hợp trong hoàn cảnh thiếu sự hướng dẫn.

## 10.3 Học tăng cường (Reinforcement learning)

- Dựa vào mẫu huấn luyện chưa biết phân lớp.
- Bộ học hành động hướng đến một hình huống và nhận kết quả phản hồi.

→ Tự giải thích phản hồi thông qua việc tạo luật và hướng dẫn.

## 11. Các mức độ học máy

- Mức độ 1: Ghi nhớ: Dự đoán mẫu dữ kiện học trước
- Mức độ 2: Lấy trung bình để xử lý nhiễu
- + Dự đoán sai lệch do thiết bị.

## 12. Phân lớp là gì?

- Phân lớp là bài toán học với một hàm mục tiêu  $f$  và



ánh xạ tập các thuộc tính  $X$  từ một nhãn lớp  $y$ .

- Một trong các thuộc tính là thuộc tính lớp.

- Hai nhãn lớp (hoặc từ lớp): Yes (1), No (0)

Input

Mô hình

Output

Tập thuộc tính  $\rightarrow$  phân lớp  $\rightarrow$  Nhãn lớp  
(X) (Y)

### 13. Học có giám sát

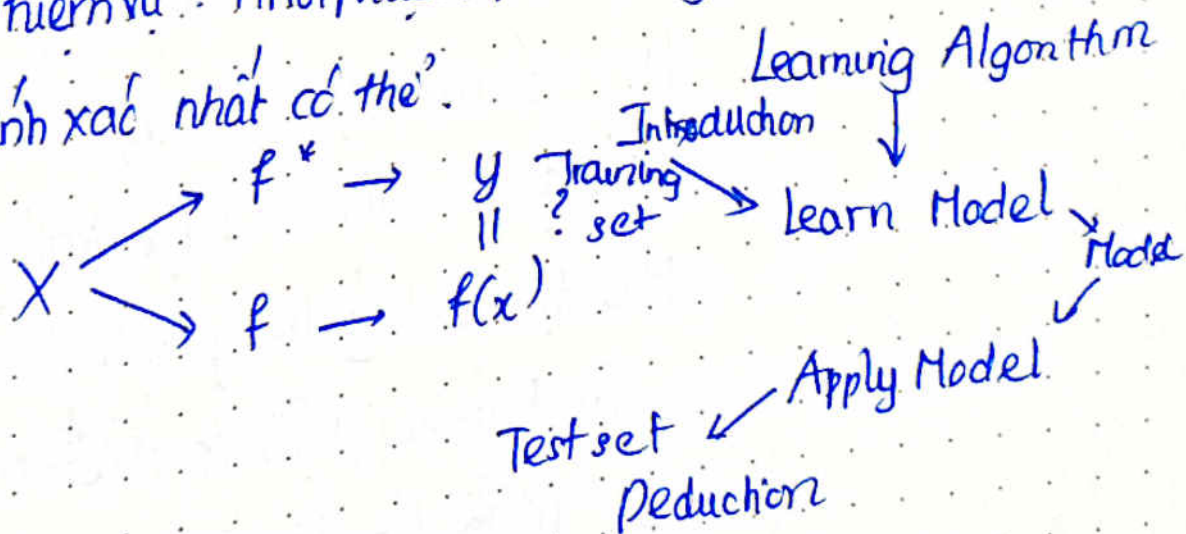
Chúng ta có các quan sát:  $x^{(1)}, x^{(2)}, \dots, x^{(N)}$  và 1 hàm  $f^*$  nào đó. Qua hàm  $f^*$  này, tạo thành các

nhãn (Label):  $y^{(1)} = f^*(x^{(1)})$   $y^{(2)} = f^*(x^{(2)})$

$\dots$   $y^{(N)} = f^*(x^{(N)})$

Nhưng hàm  $f^*$  bị mất. Chúng ta chỉ biết:  $(x^{(i)}, y^{(i)})$

Nhiệm vụ: Khôi phục lại  $f^*$  bằng mô hình  $f_w$  1 cách chính xác nhất có thể.





ints

# 14. Hunt's Algorithm

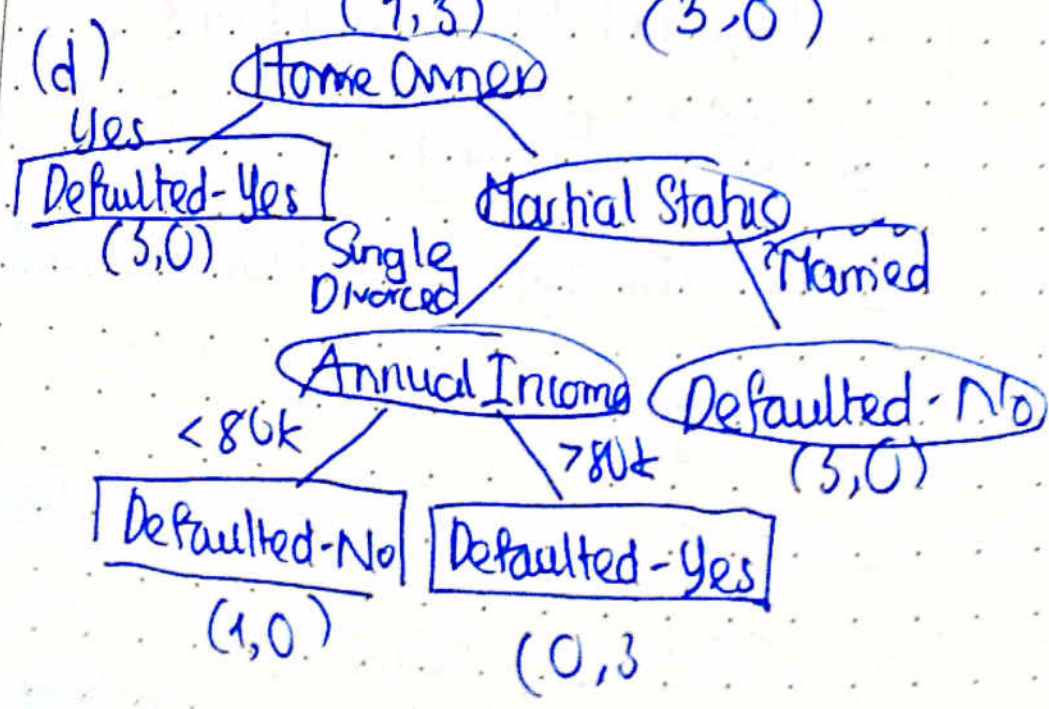
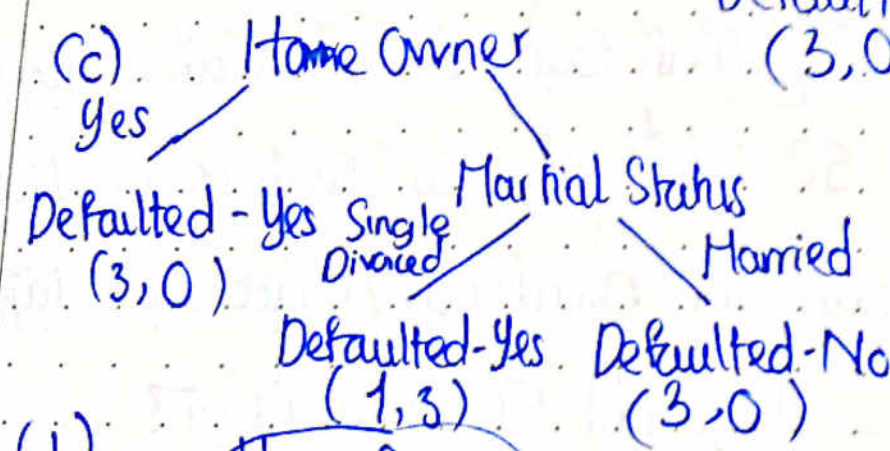
ID	Home Owner	Marital Status	Annual Income	Defaulted
1	Yes	Single	125k	No
2	No	Married	100k	No
3	No	Single	70k	No
4	Yes	Married	120k	No
5	No	Divorced	95k	No
6	No	Married	60k	Yes
7	Yes	Divorced	20k	No
8	No	No Single	85k	No
9	No	Married	75k	Yes
10	No	Single	90k	No

(a) Defaulted - No  
(7, 3)

(b) Home Owner  
Yes / No

Defaulted - Yes  
(3, 0)

Defaulted - No  
(4, 3)



## ⊗ Thuật toán ID3

- Được phát triển đồng thời bởi Quinlan & kinh nghiệm từ học nhận tạo và Breiman, Friedman, Olsh & Stone & thống kê.

- Lớp:

1. Chọn  $A \leftarrow \in$  hình quyết định "tốt nhất" cho nút kế tiếp.

2. Giá trị  $A$  là thuộc hình quyết định cho nút.

3. Với mỗi giá trị của  $A$ , tạo nhánh con mới của nút.

4. Phân loại các mẫu huấn luyện cho các nút lá.

5. Nếu các mẫu huấn luyện được phân loại hoàn toàn thì NGÚNG. Ngược lại, lặp với các nút lá mới.

- Thuộc hình tốt nhất là gì?

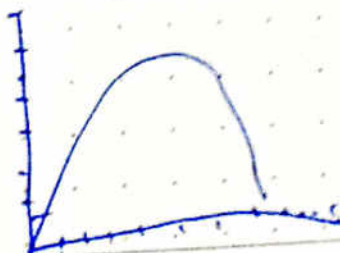
## ⊗ Độ đo Entropy

- Gọi  $S$  là tập ngẫu nhiên mẫu huấn luyện, phân tỷ lệ các mẫu  $dq$  &  $S$ .

-  $H \equiv -p \cdot \log_2 p - (1-p) \cdot \log_2 (1-p)$

- Nếu có  $n$  hơn 2 lớp?

Entropy  
(S)





Date: .....  
- Thuộc tính tốt nhất tối thiểu hoá Entropy trung bình của dữ liệu trong các nút con.

$$\text{Average Entropy}(A) = \sum_{v \in \text{value}(A)} P_v H_{A=v}$$

## 15. Cây quyết định

### 15.1. Áp dụng

Dữ liệu quan sát  $\rightarrow$  Cây quyết định  $\rightarrow$  Dự đoán - dữ liệu chưa biết  
 $\downarrow$   
Luật quyết định  $\rightarrow$

### 15.2. Thuật toán

- Các thuật toán xây dựng cây quyết định: Hunt's Algorithm, CART, ID3, C4.5, Cây định danh

### 15.3. Từ CSDL đến cây định danh

- Cơ sở dữ liệu:

- + Tập các mẫu quan sát thực tế
- + Mỗi mẫu có 1 thuộc tính cũng là thuộc tính dùng để phân loại (dự đoán / định danh đối tượng)

+ Số để huấn luyện cây quyết định

V.D: 3 màu tóc  $\times$  3 chiều cao  $\times$  3 cân nặng  $\times$  2 số liệu  
 $= 54$  tổ hợp có thể

Xác suất 1 mẫu mỗi hàng với mẫu quan sát 45%  
 Thúc đẩy: &L thuộc tính và gđ bị thuộc tính và lđ  
 → Kô thể làm = cách so sánh từ vựng

+ Hướng giải quyết

+ Đề xuất 1 thủ tục phân lớp chính xác từng mẫu

Thủ tục đúng = số lq đúng phân mẫu → c) Hết vòng  
 đúng → các mẫu mà nhãn lớp là chưa biết

VD:

Name	Hair color	Height	Weight	Lotion	Result
Sarah	blonde	average	light	no	unburned
Dana	blonde	tall	average	yes	none
Alex	brown	short	average	yes	none
Annie	blonde	short	average	no	unburned
Emily	red	short	average	no	unburned
Pete	brown	average	heavy	no	none
John	brown	average	heavy	no	none
Katie	blonde	short	light	yes	none

- Cây định danh: là g. mỗi nhánh chỉ chứa 1 kết quả

- Occam's Razor:

+ Thế giới vốn dĩ đơn giản

+ Cây định danh nhỏ nhất mà nhất quán với các mẫu

→ cây có khả năng nhất 0 và phân lớp chính xác các đtq

chưa biết: Cây định danh (1) kết hợp (2)



- Độ hỗn loạn trung bình: Average disorder

$$= \sum_b \left( \frac{n_b}{n_t} \right) \times \left( \sum_c - \frac{n_{bc}}{n_b} \log_2 \frac{n_{bc}}{n_b} \right)$$

Trong đó:

+  $n_b$  là số mẫu nhánh  $b$

+  $n_t$  là tổng số mẫu gặt cả các nhánh

+  $n_{bc}$  là tổng số mẫu  $g$  nhánh  $b \in$  lớp  $c$

- Độ hỗn loạn:

$$\text{Disorder} = \sum_c - \frac{n_{bc}}{n_b} \log_2 \frac{n_{bc}}{n_b}$$

+ Mô tả sự "hỗn loạn" ở 1 tập dữ liệu:

• DHL = 0 nếu tập dữ liệu là đồng nhất (chỉ gồm 1 lớp)

• DHL = 1 nếu tập DL chứa tất cả các lớp

- Xây dựng cây định danh:

+ Chỉ có thể xây dựng cây "nhỏ", mặc dù không đảm bảo là "nhỏ nhất"

+ Thủ tục: Chia để trị

• Tìm test cho nút gốc: Chia CSPL thành các tập con với càng nhiều mẫu  $\in$  cùng 1 lớp càng tốt.

• Với mỗi tập chứa  $n$  hơn 1 lớp, chọn 1 test  $\neq$  để chia tập

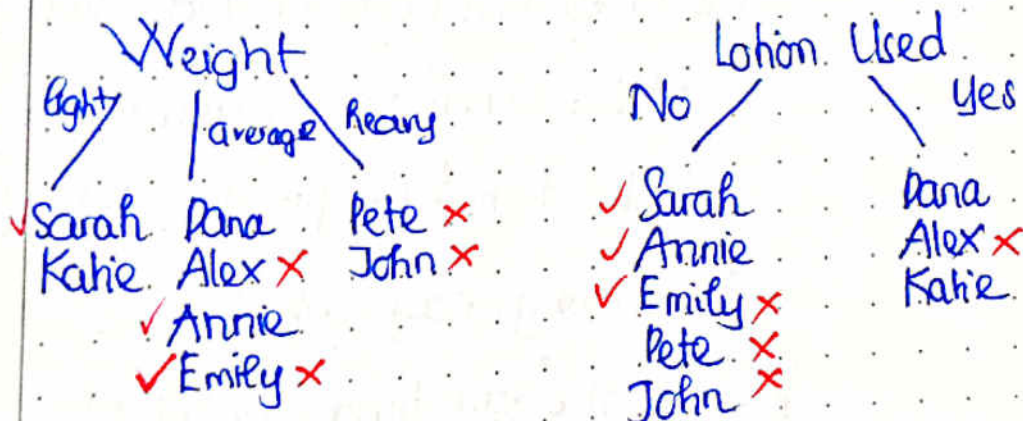
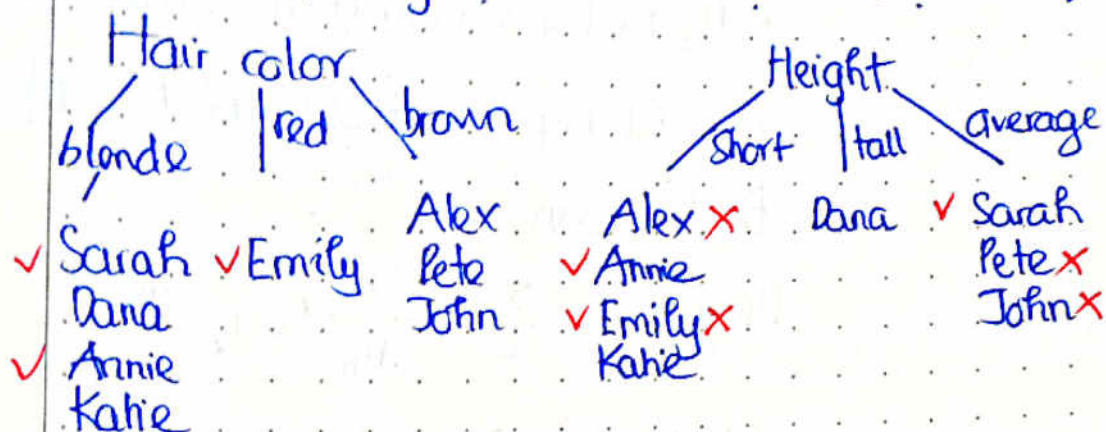
không đồng nhất  $\equiv$  các tập con đồng nhất.

+ điều kiện dừng:

• Mỗi nút lá chỉ gồm các mẫu đồng nhất.

• Không còn  $\in$  tính nào có thể phân chia nữa.

$\Rightarrow$  Thuật ngữ này cực hữu hạn độ hỗn loạn dữ liệu.



Tính độ hỗn loạn trung bình

$$\begin{aligned}
 L1: AD_{\text{Hair color}} &= \left[ \frac{4}{8} \times \left( -\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right) \right. \\
 &\quad \left. + \frac{1}{8} \times \left( -\frac{1}{1} \log_2 \frac{1}{1} - \frac{0}{4} \log_2 \frac{0}{4} \right) \right] + \frac{3}{8} \left( -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3} \right) \\
 &= 0,5 \text{ (Chon)}
 \end{aligned}$$



$$AD_{\text{Height}} = \frac{4}{8} \left( -\frac{4}{4} \log_2 \frac{2}{4} \right) + \frac{1}{8} \cdot 0 + \frac{3}{8} \left( -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right)$$

$$\approx 0,84436$$

$$AD_{\text{Weight}} = \frac{2}{8} \left( -\frac{2}{2} \log_2 \frac{1}{2} \right) + \frac{4}{8} \left( -\frac{4}{4} \log_2 \frac{2}{4} \right) + \frac{2}{8} \cdot 0$$

$$= 0,75$$

$$AD_{\text{Lotion used}} = \frac{5}{8} \left( -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) + \frac{3}{8} \cdot 0$$

$$\approx 0,6068$$

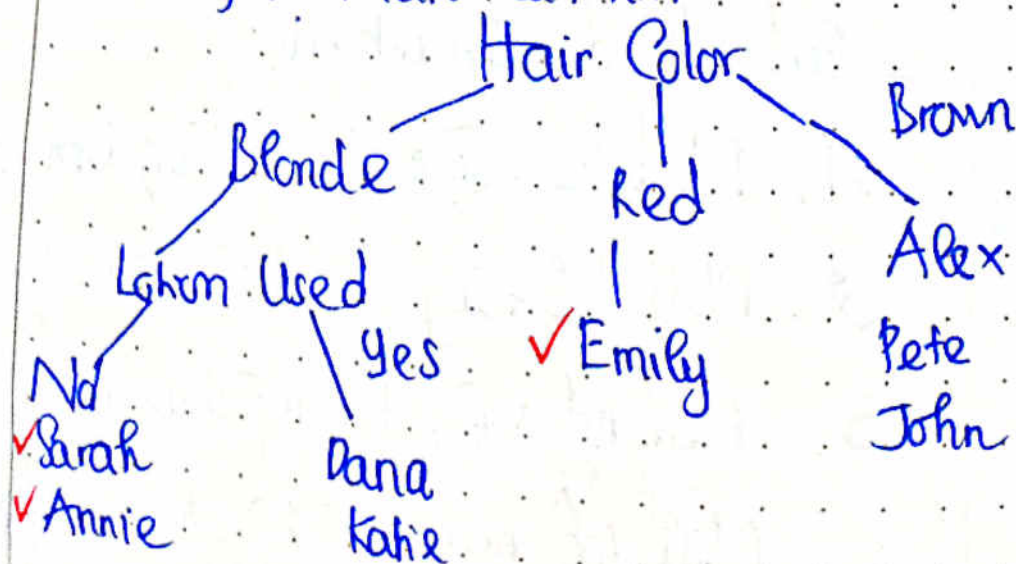
$$L2: AD_{\text{Height}} = \frac{2}{4} \left( -\frac{2}{2} \log_2 \frac{1}{2} \right) + \frac{1}{4} \cdot 0 + \frac{1}{4} \cdot 0 = 0,5$$

$$AD_{\text{Weight}} = \frac{2}{4} \left( -\frac{2}{2} \log_2 \frac{1}{2} \right) + \frac{2}{4} \left( -\frac{2}{2} \log_2 \frac{1}{2} \right) + \frac{0}{4} \cdot 0$$

$$= 1$$

$$AD_{\text{Lotion Used}} = \frac{2}{4} \cdot 0 + \frac{2}{4} \cdot 0 = 0 \text{ (Chon)}$$

Bây giờ, đánh cây này:



- Thuật hức đấm chôn + SPROUTER:

+ Lắp đèn khi: Hố nút lá chỉ gồm các mẫu đồng nhất hoặc không đồng nhất hoặc không còn ở hình nào có thể phân chia nữa.

• Chọn 1 nút lá để có tập các mẫu không đồng nhất.

• Thay thế nút lá đó bằng 1 nút test mà nó sẽ chứa tập mẫu không đồng nhất thành các tập không đồng nhất ở mức độ tối thiểu dựa vào độ đa tính hỗn loạn.

15.4 Từ cây đến luật: Rút luật từ cây, Loại bỏ các điều kiện luật không cần thiết, Loại bỏ luật thừa, Thuật hức tỉa cành (Pruner)

- Rút luật:

+ Theo dấu mũi tên dẫn từ gốc đến lá

+ Lấy các phép thử tìm kiếm

+ Các nút lá làm kết luận

1. Nếu ~~màu vàng~~ và không đúng kem thì cháy nắng

2. Nếu ~~hạt vàng~~ và đúng kem thì không cháy nắng

3. Nếu ~~không vàng~~ thì cháy nắng

4. Nếu ~~không~~ thì không cháy nắng



points

## - Thủ tục hĩa cảnh :

- + Tạo 1 luật cho mỗi đg đi từ gốc đến lá O cây quyết định
- + Đoán giá trị học 1 luật bằng cách loại bỏ ~ nên để k<sup>o</sup> ảnh hưởng KL mà cây có đc
- + Thay thế ~ luật có chung KL bằng 1 luật mặc định mà luật này sẽ đc kích hoạt khi k<sup>o</sup> có luật nào ± đc kích hoạt

## → Các tài liệu

1. Nếu k<sup>o</sup> dùng kern hoặc <sup>tốc độ thi</sup> chạy năng
2. Nếu tốc độ & dùng kern hoặc tốc độ nấu thì chạy năng