

Nhập môn Học máy và Khai phá dữ liệu (IT3190)

Nguyễn Nhật Quang

quang.nguyennhat@hust.edu.vn

Trường Đại học Bách Khoa Hà Nội
Viện Công nghệ thông tin và truyền thông
Năm học 2020-2021

Nội dung môn học:

Giới thiệu về Học máy và Khai phá dữ liệu

Tiền xử lý dữ liệu

Đánh giá hiệu năng của hệ thống

Hồi quy

Phân lớp

Cây quyết định (Decision tree)

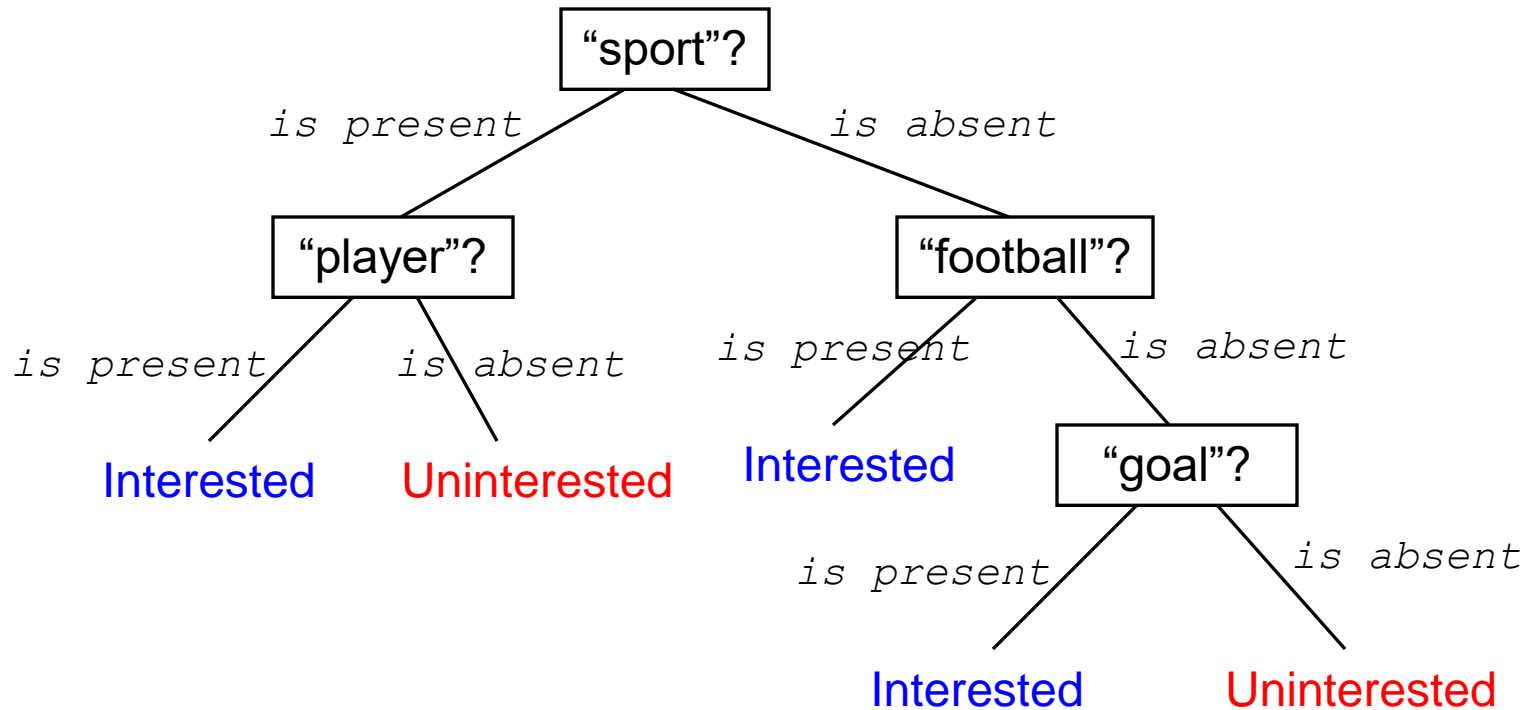
Phân cụm

Phát hiện luật kết hợp

Học cây quyết định – Giới thiệu

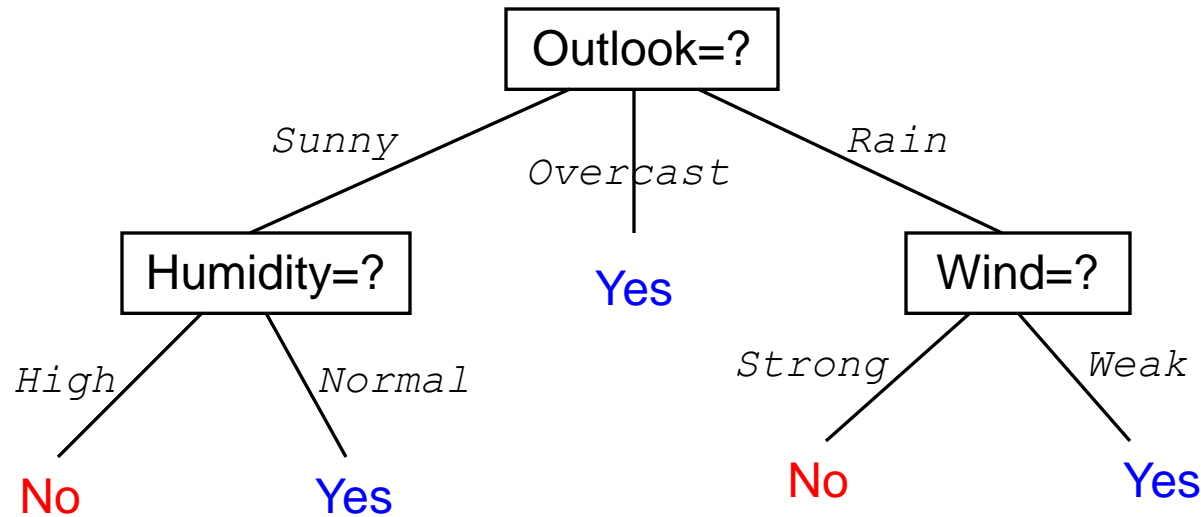
- Học cây quyết định (Decision tree –DT– learning)
 - Để học (xấp xỉ) một hàm mục tiêu có giá trị rời rạc (*discrete-valued target function*) – hàm phân lớp
 - Hàm phân lớp được biểu diễn bởi một cây quyết định
- Một cây quyết định có thể được biểu diễn (diễn giải) bằng một tập các luật IF-THEN (dễ đọc và dễ hiểu)
- Học cây quyết định có thể thực hiện ngay cả với các dữ liệu có chứa nhiễu/lỗi (noisy data)
- Là một trong các phương pháp học quy nạp (inductive learning) được dùng phổ biến nhất
- Được áp dụng thành công trong rất nhiều các bài toán ứng dụng thực tế

Ví dụ về DT: Những tin tức nào mà tôi quan tâm?



- (...,"sport",...,"player",...) → Interested
- (...,"goal",...) → Interested
- (...,"sport",...) → Uninterested

Ví dụ về DT: Một người có chơi tennis không?



- (Outlook=Overcast, Temperature=Hot, Humidity=High, Wind=Weak) → Yes
- (Outlook=Rain, Temperature=Mild, Humidity=High, Wind=Strong) → No
- (Outlook=Sunny, Temperature=Hot, Humidity=High, Wind=Strong) → No

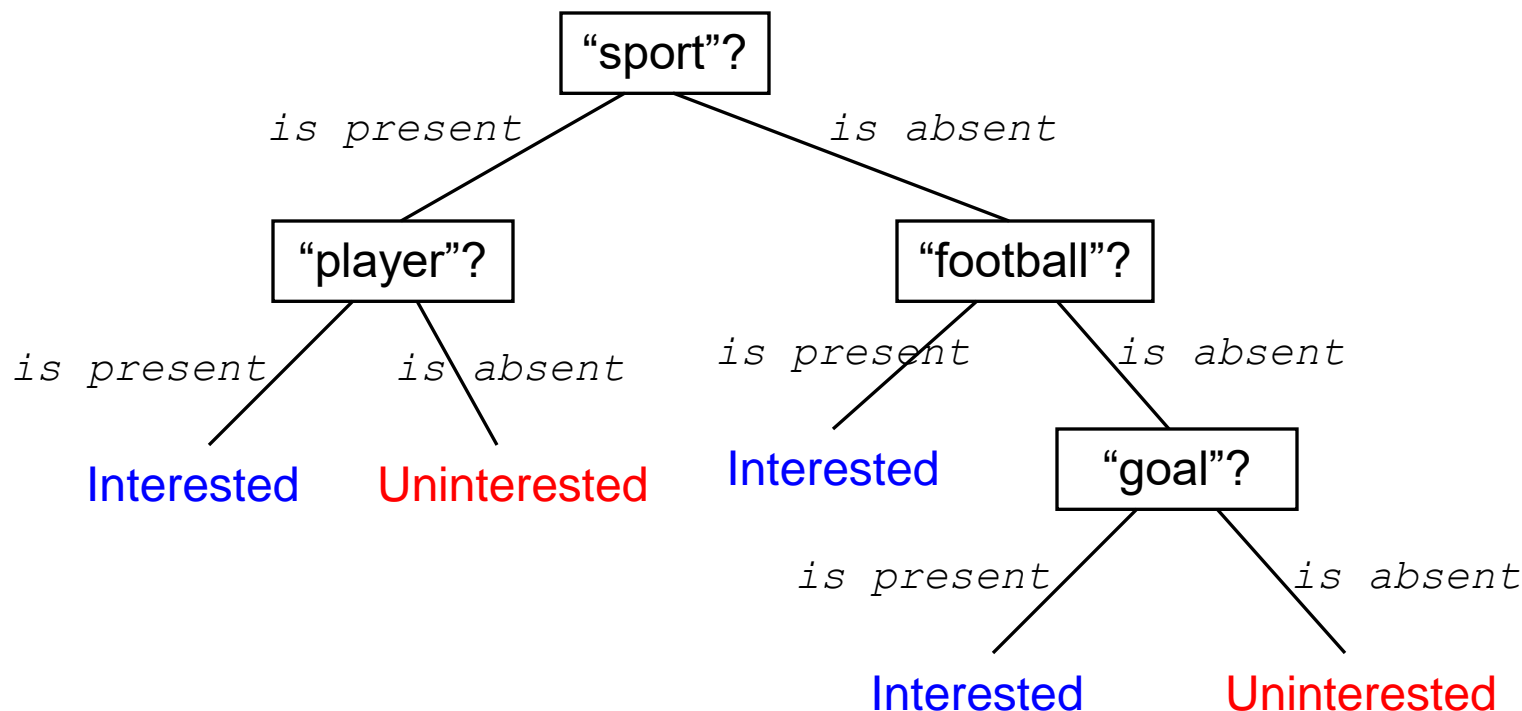
Biểu diễn cây quyết định (1)

- Mỗi nút trong (*internal node*) biểu diễn một thuộc tính cần kiểm tra giá trị (*an attribute to be tested*) đối với các ví dụ
- Mỗi nhánh (*branch*) từ một nút sẽ tương ứng với một giá trị có thể của thuộc tính gắn với nút đó
- Mỗi nút lá (*leaf node*) biểu diễn một phân lớp (*a classification*)
- Một cây quyết định học được sẽ phân lớp đối với một ví dụ, bằng cách duyệt cây từ nút gốc đến một nút lá
 - Nhãn lớp gắn với nút lá đó sẽ được gán cho ví dụ cần phân lớp

Biểu diễn cây quyết định (2)

- Một cây quyết định biểu diễn một phép tuyển (disjunction) của các kết hợp (conjunctions) của các ràng buộc đối với các giá trị thuộc tính của các ví dụ
- Mỗi đường đi (path) từ nút gốc đến một nút lá sẽ tương ứng với một kết hợp (conjunction) của các kiểm tra giá trị thuộc tính (attribute tests)
- Cây quyết định (bản thân nó) chính là một phép tuyển (disjunction) của các kết hợp (conjunctions) này
- Các ví dụ
 - Hãy xét 2 cây quyết định đã nêu ở trước...

Những tin tức nào mà tôi quan tâm?

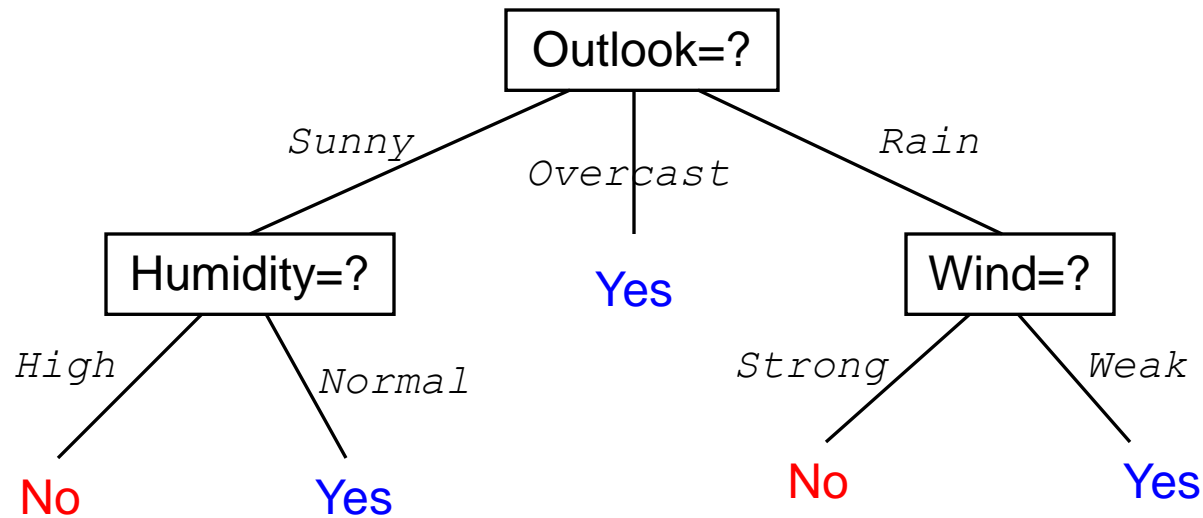


$[("sport" \text{ is present}) \wedge ("player" \text{ is present})] \vee$

$[("sport" \text{ is absent}) \wedge ("football" \text{ is present})] \vee$

$[("sport" \text{ is absent}) \wedge ("football" \text{ is absent}) \wedge ("goal" \text{ is present})]$

Một người có chơi tennis không?



$[(\text{Outlook}=\text{Sunny}) \wedge (\text{Humidity}=\text{Normal})] \vee$

$(\text{Outlook}=\text{Overcast}) \vee$

$[(\text{Outlook}=\text{Rain}) \wedge (\text{Wind}=\text{Weak})]$

Giải thuật ID3 – Ý tưởng

- Thực hiện giải thuật tìm kiếm tham lam (greedy search) đối với không gian các cây quyết định có thể
- Xây dựng (học) một cây quyết định theo chiến lược top-down, bắt đầu từ nút gốc
- Ở mỗi nút, thuộc tính kiểm tra (test attribute) là thuộc tính có khả năng phân loại tốt nhất đối với các ví dụ học gắn với nút đó
- Tạo mới một cây con (sub-tree) của nút hiện tại cho mỗi giá trị có thể của thuộc tính kiểm tra, và tập học sẽ được tách ra (thành các tập con) tương ứng với cây con vừa tạo
- Mỗi thuộc tính chỉ được phép xuất hiện tối đa 1 lần đối với bất kỳ một đường đi nào trong cây
- Quá trình phát triển (học) cây quyết định sẽ tiếp tục cho đến khi...
 - Cây quyết định phân loại hoàn toàn (perfectly classifies) các ví dụ học, hoặc
 - Tất cả các thuộc tính đã được sử dụng

Giải thuật ID3

ID3_alg(Training_Set, Class_Labels, Attributes)

Tạo nút Root của cây quyết định

If tất cả các ví dụ của Training_Set thuộc cùng lớp c , Return Cây quyết định có nút Root được gán với (có nhãn) lớp c

If Tập thuộc tính Attributes là rỗng, Return Cây quyết định có nút Root được gán với nhãn lớp \equiv **Majority_Class_Label**(Training_Set)

$A \leftarrow$ Thuộc tính trong tập Attributes có khả năng phân loại “tốt nhất” đối với Training_Set

Thuộc tính kiểm tra cho nút Root $\leftarrow A$

For each Giá trị có thể v của thuộc tính A

 Bổ sung một nhánh cây mới dưới nút Root, tương ứng với trường hợp: “Giá trị của A là v ”

 Xác định $\text{Training_Set}_v = \{\text{ví dụ } x \mid x \subseteq \text{Training_Set}, x_A = v\}$

If (Training_Set_v là rỗng) Then

 Tạo một nút lá với nhãn lớp \equiv **Majority_Class_Label**(Training_Set)

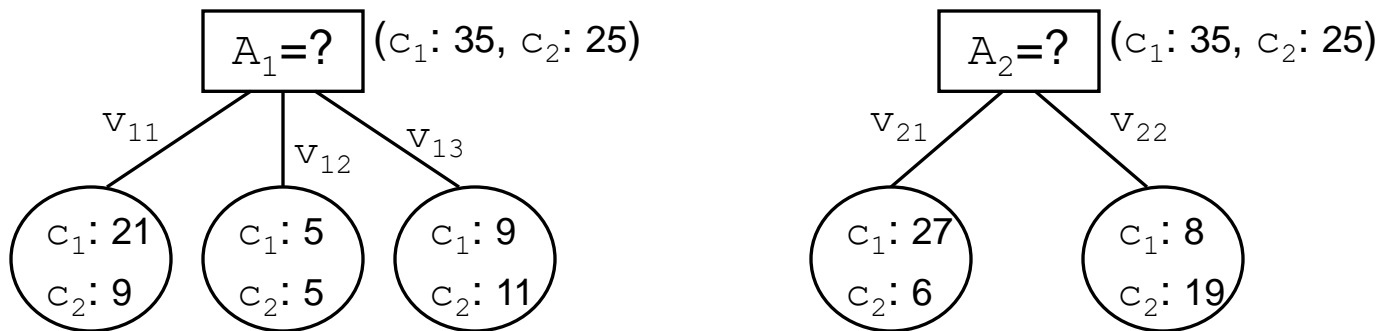
 Gắn nút lá này vào nhánh cây mới vừa tạo

Else Gắn vào nhánh cây mới vừa tạo một cây con sinh ra bởi **ID3_alg**(Training_Set $_v$, Class_Labels, {Attributes \ A })

Return Root

Lựa chọn thuộc tính kiểm tra

- Tại mỗi nút, chọn thuộc tính kiểm tra như thế nào?
- Chọn thuộc tính quan trọng nhất cho việc phân lớp các ví dụ học gắn với nút đó
- Làm thế nào để đánh giá khả năng của một thuộc tính đối với việc phân tách các ví dụ học theo nhãn lớp của chúng?
 - Sử dụng một đánh giá thống kê – *Information Gain*
- Ví dụ: Xét bài toán phân lớp có 2 lớp (c_1, c_2)
 - Thuộc tính nào, A_1 hay A_2 , nên được chọn là thuộc tính kiểm tra?



Entropy

- Một đánh giá thường được sử dụng trong lĩnh vực Information Theory
- Để đánh giá mức độ hỗn tạp (impurity/inhomogeneity) của một tập
- Entropy của tập S đối với việc phân lớp có c lớp

$$Entropy(S) = \sum_{i=1}^c -p_i \cdot \log_2 p_i$$

trong đó p_i là tỷ lệ các ví dụ trong tập S thuộc vào lớp i , và $0 \cdot \log_2 0 = 0$

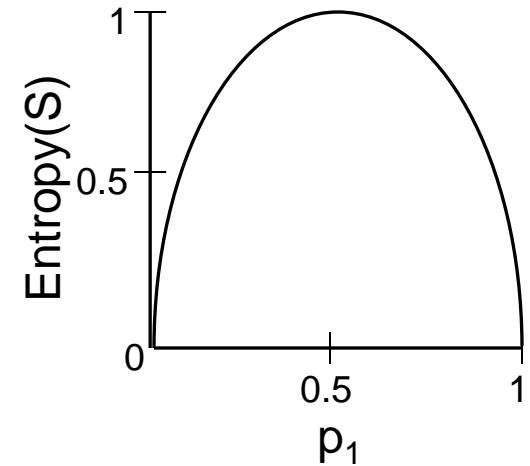
- Entropy của tập S đối với việc phân lớp có 2 lớp

$$Entropy(S) = -p_1 \cdot \log_2 p_1 - p_2 \cdot \log_2 p_2$$

- Ý nghĩa của entropy trong lĩnh vực Information Theory
 - Entropy của tập S chỉ ra số lượng bits cần thiết để mã hóa lớp của một phần tử được lấy ra ngẫu nhiên từ tập S

Entropy – Ví dụ với 2 lớp

- S gồm 14 ví dụ, trong đó 9 ví dụ thuộc về lớp c_1 và 5 ví dụ thuộc về lớp c_2
- Entropy của tập S đối với phân lớp có 2 lớp:
$$\text{Entropy}(S) = -(9/14) \cdot \log_2(9/14) - (5/14) \cdot \log_2(5/14) \approx 0.94$$
- Entropy = 0, nếu tất cả các ví dụ thuộc cùng một lớp (c_1 hoặc c_2)
- Entropy = 1, số lượng các ví dụ thuộc về lớp c_1 bằng số lượng các ví dụ thuộc về lớp c_2
- Entropy = một giá trị trong khoảng (0,1), nếu như số lượng các ví dụ thuộc về lớp c_1 khác với số lượng các ví dụ thuộc về lớp c_2



Information gain

- Information Gain của một thuộc tính đối với một tập các ví dụ:
 - Mức độ giảm về Entropy
 - Bởi việc phân chia (partitioning) các ví dụ theo các giá trị của thuộc tính đó
- Information Gain của thuộc tính A đối với tập S

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

trong đó $Values(A)$ là tập các giá trị có thể của thuộc tính A , và

$$S_v = \{x \mid x \in S, x_A = v\}$$

- Trong công thức trên, thành phần thứ 2 thể hiện giá trị Entropy sau khi tập S được phân chia bởi các giá trị của thuộc tính A
- Ý nghĩa của $Gain(S, A)$: Số lượng bits giảm được (reduced) đối với việc mã hóa lớp của một phần tử được lấy ra ngẫu nhiên từ tập S , khi biết giá trị của thuộc tính A

Tập các ví dụ học

Xét tập dữ liệu s ghi lại những ngày mà một người chơi (không chơi) tennis:

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Information Gain – Ví dụ

- Hãy tính giá trị Information Gain của thuộc tính `Wind` đối với tập học `S`
– $\text{Gain}(S, \text{Wind})$?
- Thuộc tính `Wind` có 2 giá trị có thể: `Weak` và `Strong`
- $S = \{9 \text{ ví dụ lớp Yes và } 5 \text{ ví dụ lớp No}\}$
- $S_{\text{weak}} = \{6 \text{ ví dụ lớp Yes và } 2 \text{ ví dụ lớp No có giá trị Wind=Weak}\}$
- $S_{\text{strong}} = \{3 \text{ ví dụ lớp Yes và } 3 \text{ ví dụ lớp No có giá trị Wind=Strong}\}$

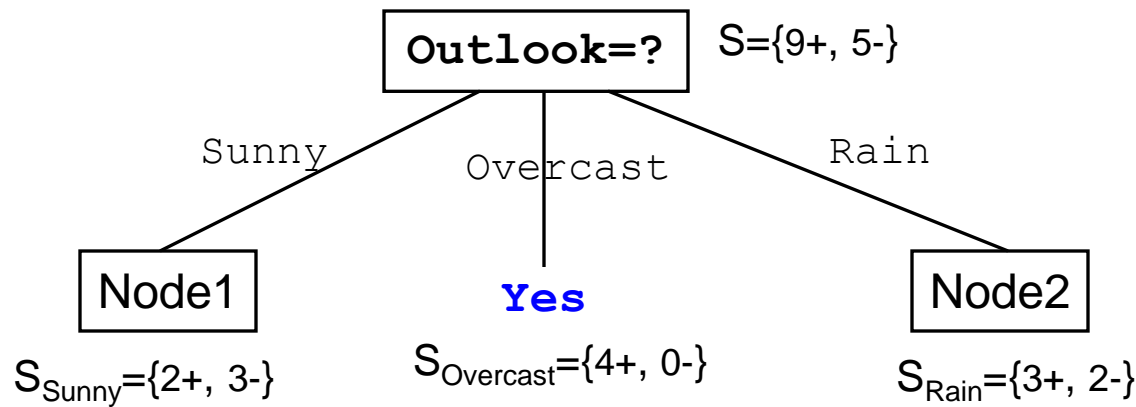
$$\begin{aligned}\text{Gain}(S, \text{Wind}) &= \text{Entropy}(S) - \sum_{v \in \{\text{Weak}, \text{Strong}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v) \\ &= \text{Entropy}(S) - (8/14) \cdot \text{Entropy}(S_{\text{Weak}}) - (6/14) \cdot \text{Entropy}(S_{\text{Strong}}) \\ &= 0.94 - (8/14) \cdot (0.81) - (6/14) \cdot (1) = 0.048\end{aligned}$$

Học cây quyết định – Ví dụ (1)

- Tại nút gốc, thuộc tính nào trong số {Outlook, Temperature, Humidity, Wind} nên được chọn là thuộc tính kiểm tra?

- $\text{Gain}(S, \text{Outlook}) = \dots = 0.246$ ← Có giá trị IG cao nhất
- $\text{Gain}(S, \text{Temperature}) = \dots = 0.029$
- $\text{Gain}(S, \text{Humidity}) = \dots = 0.151$
- $\text{Gain}(S, \text{Wind}) = \dots = 0.048$

→ Vì vậy, Outlook được chọn là thuộc tính kiểm tra cho nút gốc!



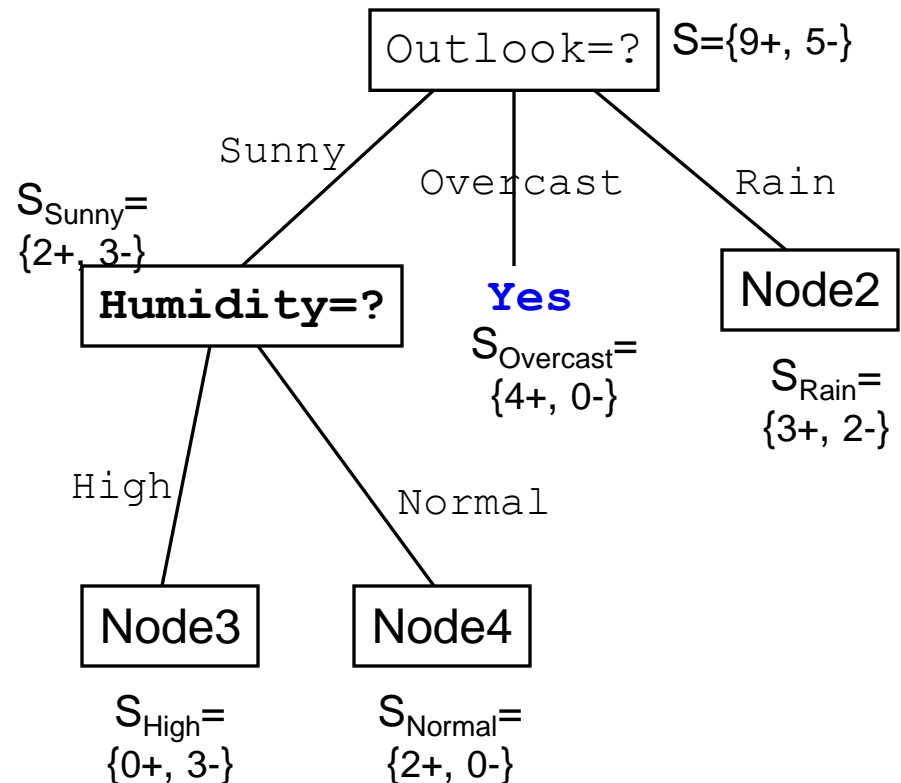
Học cây quyết định – Ví dụ (2)

- Tại nút Node1, thuộc tính nào trong số {Temperature, Humidity, Wind} nên được chọn là thuộc tính kiểm tra?

Lưu ý! Thuộc tính Outlook bị loại ra, bởi vì nó đã được sử dụng bởi cha của nút Node1 (là nút gốc)

- $\text{Gain}(S_{\text{Sunny}}, \text{Temperature}) = \dots = 0.57$
- $\text{Gain}(S_{\text{Sunny}}, \text{Humidity}) = \dots = \mathbf{0.97}$
- $\text{Gain}(S_{\text{Sunny}}, \text{Wind}) = \dots = 0.019$

→ Vì vậy, Humidity được chọn là thuộc tính kiểm tra cho nút Node1!



Học cây quyết định – Chiến lược tìm kiếm

(1)

- ID3 tìm kiếm trong không gian các giả thiết (các cây quyết định có thể) một cây quyết định phù hợp (fits) các ví dụ học
- ID3 thực hiện chiến lược tìm kiếm từ đơn giản đến phức tạp, bắt đầu với cây rỗng (empty tree)
- Quá trình tìm kiếm của ID3 được điều khiển bởi độ đo đánh giá Information Gain
- ID3 chỉ tìm kiếm một (chứ không phải tất cả các) cây quyết định phù hợp với các ví dụ học

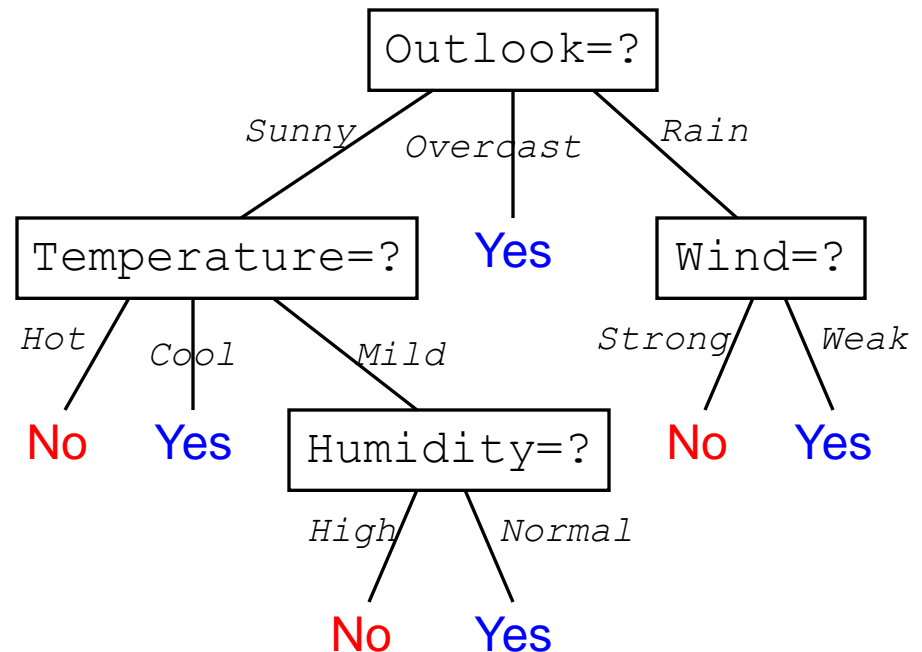
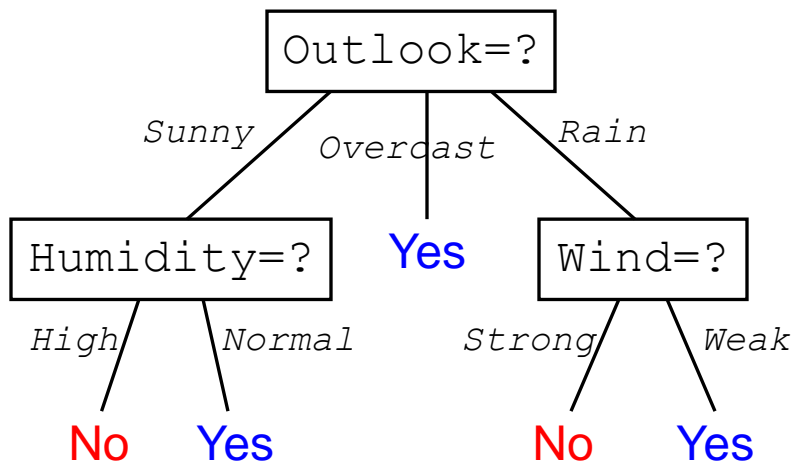
Học cây quyết định – Chiến lược tìm kiếm

(2)

- Trong quá trình tìm kiếm, ID3 không thực hiện quay lui (not backtrack)
 - Chỉ đảm bảo tìm được lời giải tối ưu cục bộ (locally optimal solution) – chứ không đảm bảo tìm được lời giải tối ưu tổng thể (globally optimal solution)
 - Một khi một thuộc tính được chọn là thuộc tính kiểm tra cho một nút, thì ID3 không bao giờ cân nhắc lại (backtracks to reconsider) lựa chọn này
- Ở mỗi bước trong quá trình tìm kiếm, ID3 sử dụng một đánh giá thống kê (Information Gain) để cải thiện giả thiết hiện tại
 - Nhờ vậy, quá trình tìm kiếm (lời giải) ít bị ảnh hưởng bởi các lỗi (nếu có) của một số ít ví dụ học

Ưu tiên trong học cây quyết định (1)

- Cả 2 cây quyết định dưới đây đều phù hợp với tập học đã cho
- Vậy thì, cây quyết định nào sẽ được ưu tiên (được học) bởi giải thuật ID3?



Ưu tiên trong học cây quyết định (2)

- Đối với một tập các ví dụ học, có thể tồn tại nhiều (hơn 1) cây quyết định phù hợp với các ví dụ học này
- Cây quyết định nào (trong số đó) được chọn?
- ID3 chọn cây quyết định phù hợp đầu tiên tìm thấy trong quá trình tìm kiếm của nó
 - Lưu ý là trong quá trình tìm kiếm, giải thuật ID3 không bao giờ cân nhắc lại các lựa chọn trước đó (without backtracking)
- Chiến lược tìm kiếm của giải thuật ID3
 - Ưu tiên các cây quyết định đơn giản (ít mức độ sâu)
 - Ưu tiên các cây quyết định trong đó một thuộc tính có giá trị *Information Gain* càng lớn thì sẽ là thuộc tính kiểm tra của một nút càng gần nút gốc

Các vấn đề trong ID3

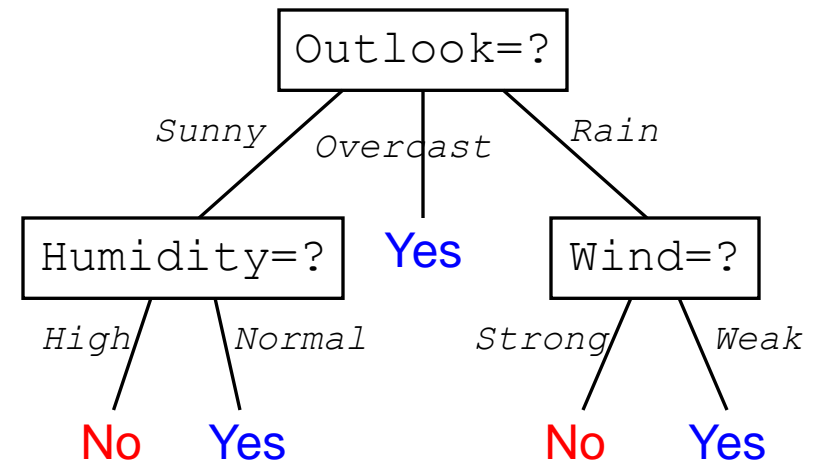
- Cây quyết định học được quá phù hợp (over-fit) với các ví dụ học
 - Xử lý các thuộc tính có kiểu giá trị liên tục (kiểu số thực)
 - Các đánh giá phù hợp hơn (tốt hơn Information Gain) đối với việc xác định thuộc tính kiểm tra cho một nút
 - Xử lý các ví dụ học thiếu giá trị thuộc tính (missing-value attributes)
 - Xử lý các thuộc tính có chi phí (cost) khác nhau
- Cải tiến của giải thuật ID3 với tất cả các vấn đề nêu trên được giải quyết: giải thuật C4.5

Over-fitting trong học cây quyết định (1)

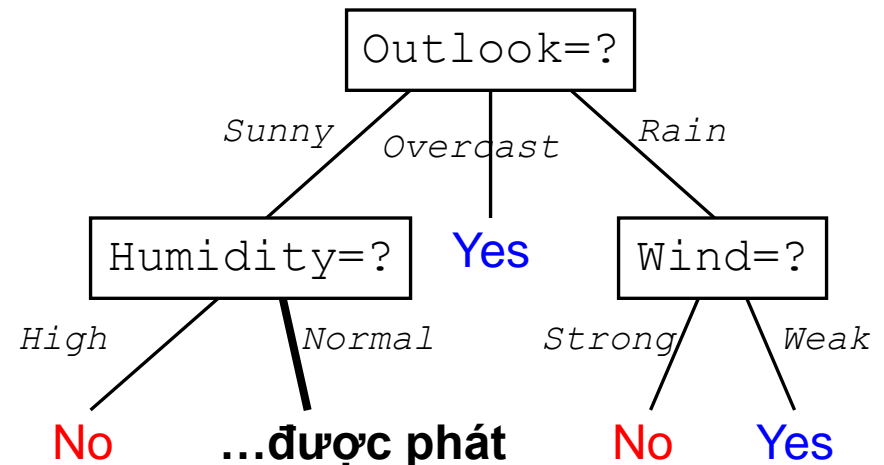
- Một cây quyết định phù hợp hoàn hảo đối với tập huấn luyện có phải là giải pháp tối ưu?
- Nếu như tập huấn luyện có nhiều/lỗi...?

Vd: Một ví dụ nhiều/lỗi (Ví dụ thực sự mang nhãn **Yes**, nhưng bị gán nhãn nhầm là **No**):

(Outlook=Sunny,
Temperature=Hot,
Humidity=Normal,
Wind=Strong, PlayTennis=**No**)



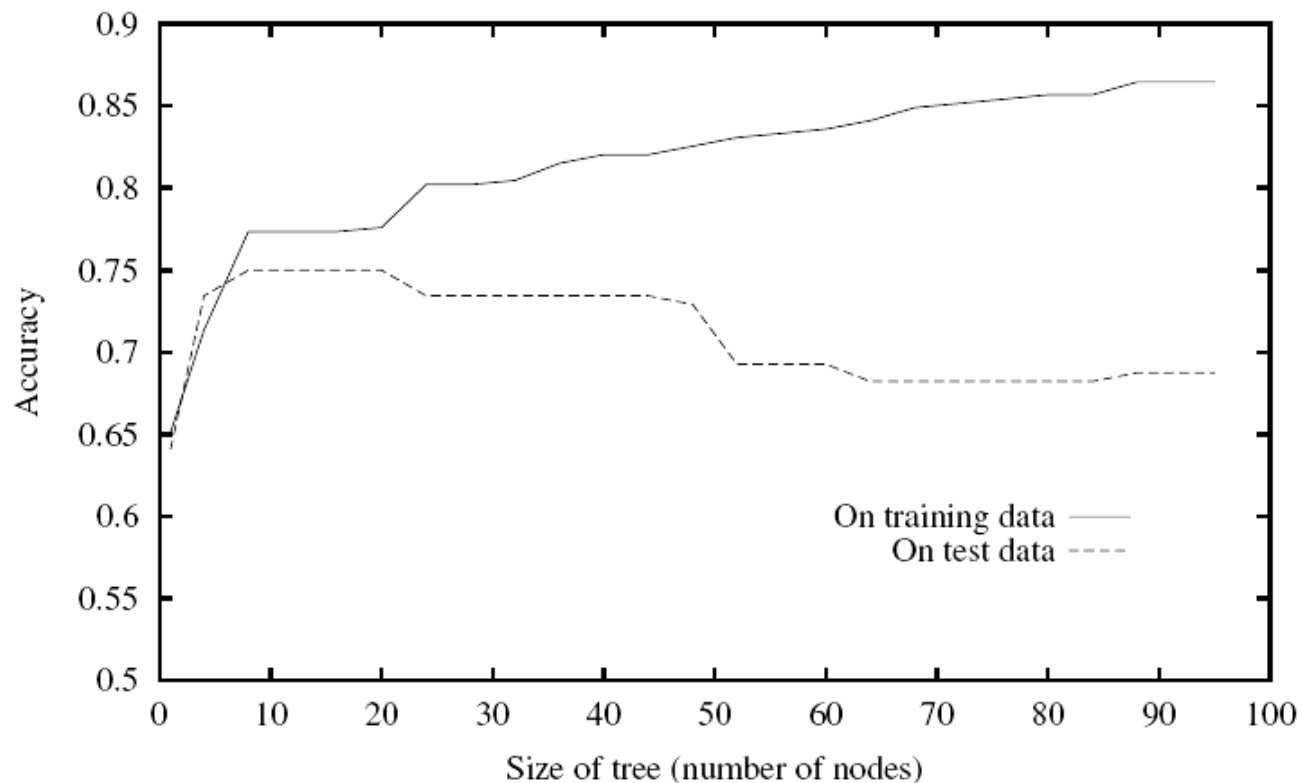
Học được một cây quyết định phức tạp hơn!
(chỉ bởi vì ví dụ nhiều/lỗi)



...được phát

Over-fitting trong học cây quyết định (1)

Tiếp tục quá trình học cây quyết định sẽ làm giảm độ chính xác đối với tập thử nghiệm mặc dù tăng độ chính xác đối với tập học



[Mitchell, 1997]

Giải quyết vấn đề over-fitting (1)

■ 2 chiến lược

- Ngừng việc học (phát triển) cây quyết định sớm hơn, trước khi nó đạt tới cấu trúc cây cho phép phân loại (khớp) hoàn hảo tập huấn luyện
- Học (phát triển) cây đầy đủ (tương ứng với cấu trúc cây hoàn toàn phù hợp đối với tập huấn luyện), và sau đó thực hiện quá trình tỉa (to post-prune) cây

■ Chiến lược tỉa cây đầy đủ (Post-pruning over-fit trees) thường cho hiệu quả tốt hơn trong thực tế

- Lý do: Chiến lược “ngừng sớm” việc học cây cần phải đánh giá chính xác được *khi nào nên ngừng việc học* (phát triển) cây – Khó xác định!

Giải quyết vấn đề over-fitting (2)

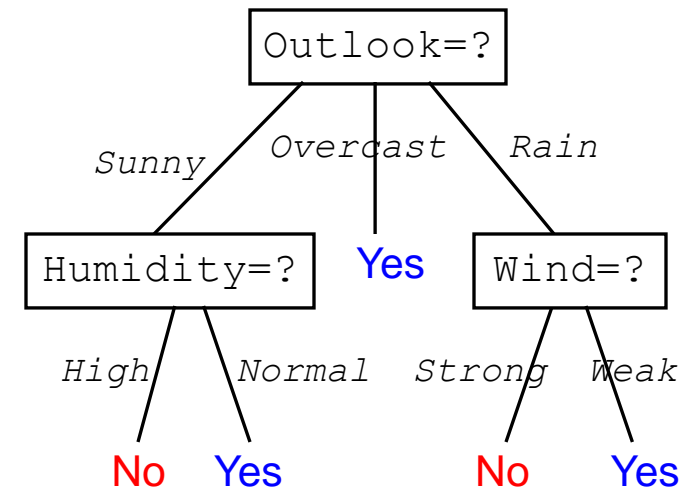
- Làm sao để chọn kích thước “phù hợp” của cây quyết định?
 - Đánh giá hiệu năng phân loại đối với một tập tối ưu (validation set)
 - Đây là phương pháp thường được sử dụng nhất
 - 2 f.f. chính: *reduced-error pruning* and *rule post-pruning*
 - Áp dụng một thí nghiệm thống kê (vd: chi-square test) để đánh giá xem việc mở rộng (hay cắt tỉa) một nút có giúp cải thiện hiệu năng đối với tập huấn luyện
 - Đánh giá độ phức tạp của việc mã hóa (thể hiện) các ví dụ học và cây quyết định, và ngừng việc học (phát triển) cây quyết định khi kích thước của việc mã hóa này là tối thiểu
 - Dựa trên nguyên lý Minimum Description Length (MDL)
 - Cần cực tiểu hóa: $\text{size}(\text{tree}) + \text{size}(\text{misclassifications}(\text{tree}))$

Reduced-error pruning

- Mỗi nút của cây (khớp hoàn toàn) được kiểm tra để cắt tỉa
- Một nút sẽ bị cắt tỉa nếu cây (sau khi cắt tỉa nút đó) đạt được hiệu năng không tồi hơn cây ban đầu đối với tập tối ưu (validation set)
- Cắt tỉa một nút bao gồm các việc:
 - Loại bỏ toàn bộ cây con (sub-tree) gắn với nút bị cắt tỉa
 - Chuyển nút bị cắt tỉa thành một nút lá (nhãn phân lớp)
 - Gắn với nút lá này (nút bị cắt tỉa) nhãn lớp chiếm số đông trong tập huấn luyện gắn với nút đó
- Lặp lại việc cắt tỉa các nút
 - Luôn lựa chọn một nút mà việc cắt tỉa nút đó tối đa hóa khả năng phân loại của cây quyết định đối với tập tối ưu (validation set)
 - Kết thúc, khi việc cắt tỉa thêm nút làm giảm khả năng phân loại của cây quyết định đối với tập tối ưu (validation set)

Rule post-pruning

- Học (phát triển) cây quyết định hoàn toàn phù hợp với tập huấn luyện
- Chuyển biểu diễn cây quyết định học được thành một tập các luật tương ứng (tạo một luật cho mỗi đường đi từ nút gốc đến một nút lá)
- Rút gọn (tổng quát hóa) mỗi luật (độc lập với các luật khác), bằng cách loại bỏ bất kỳ điều kiện nào giúp mang lại sự cải thiện về hiệu quả phân loại của luật đó
- Sắp xếp các luật đã rút gọn theo khả năng (hiệu quả) phân loại, và sử dụng thứ tự này cho việc phân loại các ví dụ trong tương lai



IF (Outlook=Sunny) \wedge
(Humidity=Normal)
THEN (PlayTennis=Yes)

Các thuộc tính có giá trị liên tục

- Cần xác định (chuyển đổi thành) các thuộc tính có giá trị rời rạc, bằng cách chia khoảng giá trị liên tục thành một tập các khoảng (intervals) không giao nhau
- Đối với thuộc tính (có giá trị liên tục) A , tạo một thuộc tính mới kiểu nhị phân A_v sao cho: A_v là đúng nếu $A > v$, và là sai nếu ngược lại
- Làm thế nào để xác định giá trị ngưỡng v “tốt nhất”?
 - Chọn giá trị ngưỡng v giúp sinh ra giá trị *Information Gain* cao nhất
- Ví dụ:
 - Sắp xếp các ví dụ học theo giá trị tăng dần đối với thuộc tính *Temperature*
 - Xác định các ví dụ học liên tiếp nhưng khác phân lớp
 - Có 2 giá trị ngưỡng có thể: Temperature_{54} và Temperature_{85}
 - Thuộc tính mới kiểu nhị phân Temperature_{54} được chọn, bởi vì $\text{Gain}(S, \text{Temperature}_{54}) > \text{Gain}(S, \text{Temperature}_{85})$

Temperature	40	48	60	72	80	90
PlayTennis	No	No	Yes	Yes	Yes	No

Các đánh giá khác cho lựa chọn thuộc tính

■ Xu hướng của đánh giá *Information Gain*

→ Ưu tiên các thuộc tính có nhiều giá trị hơn các thuộc tính có ít giá trị

Vd: Thuộc tính `Date` có số lượng rất lớn các giá trị có thể

- Thuộc tính này sẽ có giá trị *Information Gain* cao nhất
- Một mình thuộc tính này có thể phân loại hoàn hảo toàn bộ tập huấn luyện (thuộc tính này phân chia các ví dụ học thành rất nhiều các tập con có kích thước rất nhỏ)
- Thuộc tính này được chọn là thuộc tính kiểm tra ở nút gốc (của cây quyết định chỉ có mức độ sâu bằng 1, nhưng rất rộng, rất nhiều phân nhánh)

■ Một đánh giá khác: *Gain Ratio*

→ Giảm ảnh hưởng của các thuộc tính có (rất) nhiều giá trị

$$GainRatio(S, A) = \frac{Gain(S, A)}{SplitInformation(S, A)}$$

$$SplitInformation(S, A) = - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} \log_2 \frac{|S_v|}{|S|}$$

(trong đó $Values(A)$ là tập các giá trị có thể của thuộc tính A , và $S_v = \{x \mid x \in S, x_A = v\}$)

Xử lý các thuộc tính thiếu giá trị (1)

- Giả sử thuộc tính A là một ứng cử cho thuộc tính kiểm tra ở nút n
 - Xử lý thế nào với ví dụ x không có (thiếu) giá trị đối với thuộc tính A (tức là: x_A là không xác định)?
 - Gọi S_n là tập các ví dụ học gắn với nút n có giá trị đối với thuộc tính A
- Giải pháp 1: x_A là giá trị phổ biến nhất đối với thuộc tính A trong số các ví dụ thuộc tập S_n
- Giải pháp 2: x_A là giá trị phổ biến nhất đối với thuộc tính A trong số các ví dụ thuộc tập S_n có cùng phân lớp với x

Xử lý các thuộc tính thiếu giá trị (2)

→ Giải pháp 3:

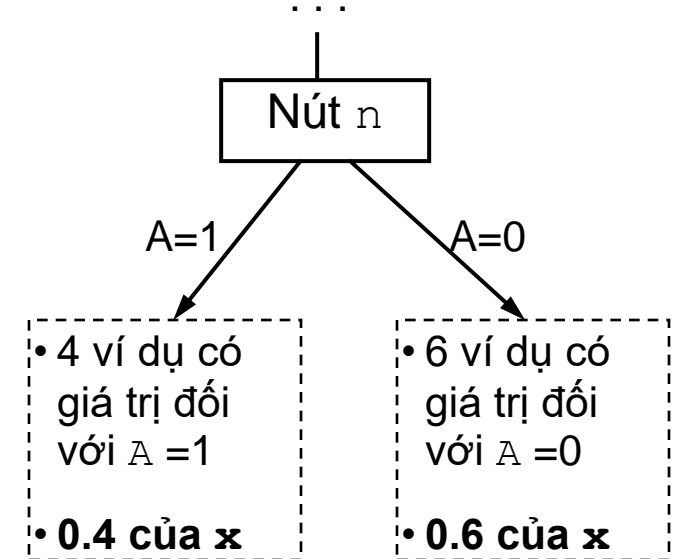
- Tính xác suất p_v đối với mỗi giá trị có thể v của thuộc tính A
- Gán *phần (fraction)* p_v của ví dụ x đối với nhánh tương ứng của nút n
- *Những ví dụ một phần (fractional instances)* này được sử dụng để tính giá trị *Information Gain*

Ví dụ:

- Một thuộc tính kiểu nhị phân (0/1) A
- Nút n bao gồm:
 - Một ví dụ x (thiếu giá trị đối với A)
 - 4 ví dụ có giá trị đối với A bằng 1, và
 - 6 ví dụ có giá trị đối với A bằng 0

$$p(x_A=1) = 4/10 = 0.4$$

$$p(x_A=0) = 6/10 = 0.6$$



Các thuộc tính có chi phí khác nhau

- Trong một số bài toán học máy, các thuộc tính có thể được gắn với các chi phí (độ quan trọng) khác nhau
 - Ví dụ: Trong việc học để phân loại các bệnh y tế, `BloodTest` có chi phí \$150, trong khi `TemperatureTest` có chi phí \$10
- Xu hướng học các cây quyết định
 - Sử dụng càng nhiều càng tốt các thuộc tính có chi phí thấp
 - Chỉ sử dụng các thuộc tính có chi phí cao khi cần thiết (để giúp đạt được các phân loại đáng tin cậy)
- Làm sao để học một cây quyết định với chi phí thấp?
 - Sử dụng các đánh giá khác IG cho việc xác định thuộc tính kiểm tra

$$\frac{Gain^2(S, A)}{Cost(A)}$$

[Tan and Schlimmer, 1990]

$$\frac{2^{Gain(S, A)} - 1}{(Cost(A) + 1)^w}$$

[Nunez, 1988; 1991]

($w \in [0, 1]$) là hằng số xác định mức độ quan trọng giữa chi phí và *Information Gain*)

Học cây quyết định – Khi nào?

- Các ví dụ học được biểu diễn bằng các cặp (thuộc tính, giá trị)
 - Phù hợp với các thuộc tính có giá trị rời rạc
 - Đối với các thuộc tính có giá trị liên tục, phải rời rạc hóa
- Hàm mục tiêu có giá trị đầu ra là các giá trị rời rạc
 - Ví dụ: Phân loại các ví dụ vào lớp phù hợp
- Rất phù hợp khi hàm mục tiêu được biểu diễn ở dạng tuyển hoặc (disjunctive form)
- Tập huấn luyện có thể chứa nhiều/lỗi
 - Lỗi trong phân loại (nhãn lớp) của các ví dụ học
 - Lỗi trong giá trị thuộc tính biểu diễn các ví dụ học
- Tập huấn luyện có thể chứa các thuộc tính thiếu giá trị
 - Các giá trị đối với một thuộc tính là không xác định đối với một số ví dụ học

Tài liệu tham khảo

- T. M. Mitchell. *Machine Learning*. McGraw-Hill, 1997.
- M. Nunez. *Economic induction: A case study*. In Proceedings of the 3rd European Working Session on Learning, EWSL-88, pp.139-145. California: Morgan Kaufmann, 1988.
- M. Nunez. *The use of background knowledge in decision tree induction*. Machine Learning, 6(3): 231-250, 1991.
- M. Tan and J. C. Schlimmer. *Two case studies in cost-sensitive concept acquisition*. In Proceedings of the 8th National Conference on Artificial Intelligence, AAAI-90, pp.854-860, 1990.