

NHẬP MÔN TRÍ TUỆ NHÂN TẠO

HỌC MÁY

ThS. Vũ Hoài Thư

Ngày 19 tháng 4 năm 2024



Nội dung

- 1 Giới thiệu
- 2 Học cây quyết định
- 3 Phân loại Bayes đơn giản
- 4 Học dựa trên ví dụ

Giới thiệu

Tài liệu tham khảo

- N. Nilsson. Introduction to machine learning:
<http://ai.stanford.edu/people/nilsson/mlbook.html>
- T. Mitchell. Machine learning. McGraw-Hill, 1997.
- E. Alpaydin. Introduction to machine learning. MIT Press, 2004.
- M. Mohri, A. Rostamizadeh, A. Talwalkar. Foundations of Machine Learning. MIT Press, 2012
- Vũ Hữu Tiệp. Machine Learning cơ bản
<https://machinelearningcoban.com/about/>

Công cụ và dữ liệu

- Bộ công cụ Weka: <http://www.cs.waikato.ac.nz/ml/weka>
- Kho dữ liệu mẫu UC Irvine:
<http://www.ics.uci.edu/mlearn/ML/Repository.html>

Một số ứng dụng của học máy(1/3)

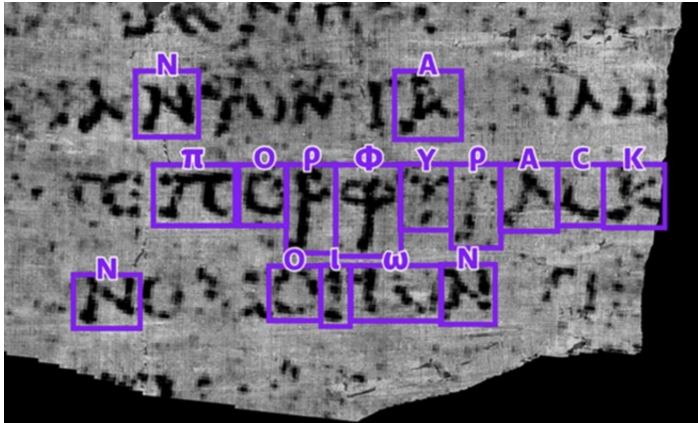
- Những ứng dụng khó lập trình theo cách thông thường do không tồn tại hoặc khó giải thích kinh nghiệm, kỹ năng của con người
 - Nhận dạng chữ viết, âm thanh, hình ảnh
 - Lái xe tự động, thám hiểm sao Hỏa
- Chương trình máy tính có khả năng thích nghi: lời giải thay đổi theo thời gian hoặc theo tình huống cụ thể
 - Chương trình trợ giúp cá nhân
 - Định tuyến mạng

Một số ứng dụng của học máy(2/3)

- Khai phá (phân tích) dữ liệu
 - Hồ sơ bệnh án → tri thức y học
 - Dữ liệu bán hàng → Quy luật kinh doanh

Một số ứng dụng của học máy(3/3)

Hầu hết các ứng dụng trí tuệ nhân tạo ngày nay có sử dụng học máy



Hình: Trí tuệ nhân tạo giải mã văn tự trên 2000 năm tuổi

Học máy là gì?

- Học máy (Machine Learning–ML) là một tập con của trí tuệ nhân tạo.
- Nó là một lĩnh vực nhỏ trong khoa học máy tính, có khả năng tự học hỏi dựa trên dữ liệu được đưa vào mà không cần phải được lập trình cụ thể
- Học máy là khả năng của chương trình máy tính sử dụng kinh nghiệm, quan sát, hoặc dữ liệu trong quá khứ để cải thiện các công việc trong tương lai thay vì chỉ thực hiện theo đúng các quy tắc đã được lập trình sẵn.

Ví dụ

- Học nhận dạng chữ:
 - Task (T): Nhận dạng chữ cái từ hình ảnh
 - Performance (P): Phần trăm chữ nhận dạng đúng
 - Experience (E): Ảnh số của chữ và chữ tương ứng
- Dịch máy:
 - Task (T): Dịch một câu tiếng Anh sang tiếng Việt
 - Performance (P): Độ đo dịch máy (ví dụ số câu đúng, số mệnh đề đúng,...)
 - Experience (E): Cặp câu tiếng Anh và tiếng Việt tương ứng

Vấn đề cần quan tâm (1/3)

- Kinh nghiệm hoặc dữ liệu cho học máy được cho dưới dạng nào?
- Lựa chọn biểu diễn cho hàm đích ra sao?

Vấn đề cần quan tâm (2/3)

Việc sử dụng những dạng kinh nghiệm và dạng biểu diễn khác nhau dẫn tới những dạng học máy khác nhau:

- Học có giám sát (supervised learning)
- Học không giám sát (un-supervised learning)
- Học bán giám sát (semi-supervised learning)
- Học tăng cường (reinforcement learning)

Vấn đề cần quan tâm (3/3)

Lựa chọn biểu diễn cho hàm đích ra sao?

- Sử dụng hàm: $y = w_1x_1 + w_2x_2 + \dots + w_nx_n$
- Sử dụng các luật
- Sử dụng mạng nơron
- Sử dụng cây quyết định
- Sử dụng các mô hình xác suất

Một số khái niệm

- **Mẫu** hay ví dụ (samples): là đối tượng cần xử lý (ví dụ phân loại)
 - Ví dụ: Khi lọc thư rác thì mỗi thư là một mẫu
- Mẫu thường được mô tả bằng tập thuộc tính hay **đặc trưng** (features)
 - Ví dụ: trong chuẩn đoán bệnh, thuộc tính là triệu chứng của người bệnh, và các tham số khác như chiều cao, cân nặng,...
- **Nhãn** phân loại (label): Thể hiện loại của đối tượng mà ta cần dự đoán
 - Ví dụ: nhãn phân loại thư rác có thể là “rác” hoặc “bình thường”

Ví dụ

thuộc tính

nhãn

mẫu

| Ngày | Trời | Nhiệt độ | Độ ẩm | Gió | Chơi tennis |
|------|------|------------|-------------|------|-------------|
| D1 | nắng | nóng | cao | yếu | không |
| D2 | nắng | nóng | cao | mạnh | không |
| D3 | u ám | nóng | cao | yếu | có |
| D4 | mưa | trung bình | cao | yếu | có |
| D5 | mưa | lạnh | bình thường | yếu | có |
| D6 | mưa | lạnh | bình thường | mạnh | không |
| D7 | u ám | lạnh | bình thường | mạnh | có |
| D8 | nắng | trung bình | cao | yếu | không |
| D9 | nắng | lạnh | bình thường | yếu | có |
| D10 | mưa | trung bình | bình thường | yếu | có |
| D11 | nắng | trung bình | bình thường | mạnh | có |
| D12 | u ám | trung bình | cao | mạnh | có |
| D13 | u ám | nóng | bình thường | yếu | có |
| D14 | mưa | trung bình | cao | mạnh | không |

Một số dạng học máy phổ biến

- Học có giám sát (supervised learning): Là dạng học máy trong đó cho trước tập dữ liệu huấn luyện dưới dạng các ví dụ cùng với giá trị đầu ra hay giá trị đích.
 - Phân loại (classification)
 - Hồi quy (regression)
- Học không giám sát (un-supervised learning): Là dạng học máy trong đó các ví dụ được cung cấp nhưng không có giá trị đầu ra hay giá trị đích.
 - Học luật kết hợp (association)
 - Phân cụm (clustering)
- Học bán giám sát (semi-supervised learning): Là dạng học máy trong đó chỉ một phần tập dữ liệu huấn luyện dưới dạng các ví dụ được cho cùng với giá trị đầu ra hay giá trị đích.
- Học tăng cường (reinforcement learning)

Phân loại (Classification)

Giá trị đích là các giá trị rời rạc

| Dung tích động cơ | Loại xe | Phân khúc |
|-------------------|---------|------------|
| 3200 | Sedan | Cao cấp |
| 2500 | Sedan | Cao cấp |
| 2500 | SUV | Trung bình |
| 2000 | Sedan | Trung bình |
| 3500 | SUV | Cao cấp |
| 1800 | Sedan | Trung bình |

Hồi quy (Regression)

Giá trị đích là các giá trị liên tục

| Dung tích động cơ | Tuổi của xe | Giá bán (triệu đồng) |
|-------------------|-------------|----------------------|
| 3200 | 1 | 2500 |
| 2500 | 2 | 1600 |
| 2500 | 4 | 1300 |
| 2000 | 5 | 600 |
| 1800 | 1 | 915 |
| 1800 | 3 | 725 |

Ứng dụng: Dự đoán giá xe, giá vàng, chứng khoán,...

Học luật kết hợp (Association)

- Ví dụ: Phân tích giao dịch, mua bán (hoá đơn mua hàng)
- $P(Y|X)$: Xác suất người mua hàng X còn mua hàng Y
- Ví dụ luật kết hợp
 - Người mua bánh mì thường mua bơ
 - Người mua lạc rang thường mua bia

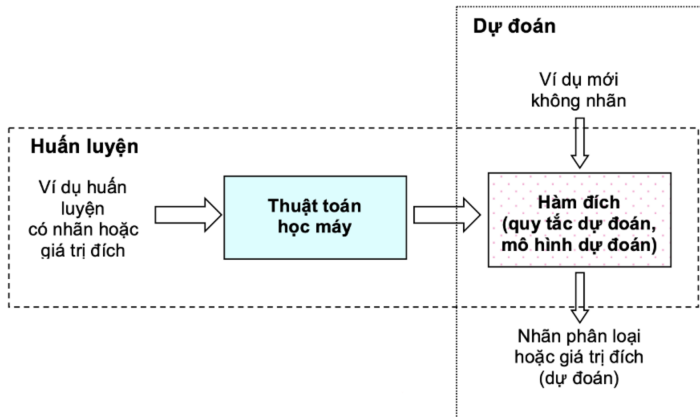
Phân cụm (Clustering)

- Nhóm những trường hợp tương tự với nhau
- Không có giá trị đầu ra
- Ứng dụng:
 - Phân cụm khách hàng, phân cụm sinh viên
 - Phân đoạn ảnh
 - Thiết kế vi mạch

Học tăng cường (Reinforcement learning)

- Kinh nghiệm không được cho trực tiếp dưới dạng đầu vào / đầu ra
- Hệ thống nhận được một giá trị thưởng (reward) là kết quả cho một chuỗi hành động nào đó
- Thuật toán cần học cách hành động để cực đại hóa giá trị thưởng

Hệ thống học máy điển hình



Học cây quyết định

Giới thiệu

- Học cây quyết định (Decision Tree) được sử dụng để học một hàm mục tiêu có giá trị rời rạc
- Hàm phân lớp được biểu diễn bởi một cây quyết định
- Một cây quyết định có thể biểu diễn (diễn giải) bằng 1 tập các luật IF-THEN
- Học cây quyết định có thể thực hiện ngay cả với các dữ liệu có chứa nhiễu/lỗi
- 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

Dữ liệu huấn luyện

thuộc tính

nhãn

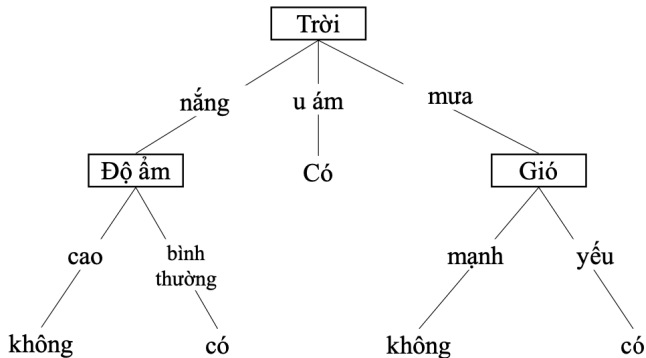
mẫu

| Ngày | Trời | Nhiệt độ | Độ ẩm | Gió | Chơi tennis |
|------|------|------------|-------------|------|-------------|
| D1 | nắng | nóng | cao | yếu | không |
| D2 | nắng | nóng | cao | mạnh | không |
| D3 | u ám | nóng | cao | yếu | có |
| D4 | mưa | trung bình | cao | yếu | có |
| D5 | mưa | lạnh | bình thường | yếu | có |
| D6 | mưa | lạnh | bình thường | mạnh | không |
| D7 | u ám | lạnh | bình thường | mạnh | có |
| D8 | nắng | trung bình | cao | yếu | không |
| D9 | nắng | lạnh | bình thường | yếu | có |
| D10 | mưa | trung bình | bình thường | yếu | có |
| D11 | nắng | trung bình | bình thường | mạnh | có |
| D12 | u ám | trung bình | cao | mạnh | có |
| D13 | u ám | nóng | bình thường | yếu | có |
| D14 | mưa | trung bình | cao | mạnh | không |

Dữ liệu

- n mẫu huấn luyện, mỗi mẫu là một cặp $\langle \mathbf{x}, y \rangle$
 - \mathbf{x} là vector thuộc tính
 - y là nhãn phân loại, $y \in C$ (tập các nhãn)
- Ví dụ mẫu D4
 - $\mathbf{x} = (\text{mưa}, \text{trung bình}, \text{cao}, \text{yếu})$
 - $y = \text{có}$

Ví dụ cây quyết định



Biểu diễn cây quyết định

- Mỗi nút trong (internal node) biểu diễn một thuộc tính cần kiểm tra giá trị 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.
- Một cây quyết định học được sẽ được phân lớp đối với một ví dụ, bằng cách duyệt từ nút gốc tới nút lá
→ Nhãn lớp gắn với nút lá đó sẽ được gán cho ví dụ cần phân lớp.

Biểu diễn dưới dạng quy tắc

- Cây quyết định có thể biểu diễn tương đương dưới dạng các quy tắc logic
- Mỗi cây là tuyển của các quy tắc, mỗi quy tắc bao gồm các phép hội
- Ví dụ:

$(\text{Trời} = \text{nắng} \wedge \text{Độ ẩm} = \text{bình_thường})$
 $\vee (\text{Trời} = \text{u_ám})$
 $\vee (\text{Trời} = \text{mưa} \wedge \text{Gió} = \text{yếu})$

Học cây quyết định

- Cây quyết định được học (xây dựng) từ dữ liệu huấn luyện
- Với mỗi bộ dữ liệu có thể xây dựng nhiều cây quyết định
 - Chọn cây nào?

Quá trình học là quá trình tìm kiếm cây quyết định phù hợp với dữ liệu huấn luyện

- Cho phép phân loại đúng dữ liệu huấn luyện

Thuật toán 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 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 kì một đường đi nào trong cây**

Thuật toán ID3

- Xây dựng lần lượt các nút của cây bắt đầu từ gốc
- Thuật toán
 - **Khởi đầu:** nút hiện thời là nút gốc chứa toàn bộ tập dữ liệu huấn luyện
 - Tại nút hiện thời n , lựa chọn thuộc tính:
 - Chưa được sử dụng ở nút tổ tiên
 - Cho phép phân chia tập dữ liệu hiện thời thành các tập con **một cách tốt nhất**
 - Với mỗi giá trị thuộc tính được chọn thêm một nút con bên dưới
 - Chia các ví dụ ở nút hiện thời về các nút con theo giá trị thuộc tính được chọn
 - **Lặp** (đệ quy) cho tới khi:
 - Tất cả các thuộc tính đã được sử dụng ở các nút phía trên, hoặc
 - Tất cả ví dụ tại nút hiện thời có cùng nhãn phân loại
 - Nhãn của nút được lấy theo đa số nhãn của ví dụ tại nút hiện thời

Vấn đề: Lựa chọn thuộc tính tại mỗi nút như thế nào?

Tiêu chuẩn chọn thuộc tính của ID3

- Tại mỗi nút n
 - Tập (con) dữ liệu ứng với nút đó
 - Cần lựa chọn thuộc tính cho phép phân chia tập dữ liệu tốt nhất
- Tiêu chuẩn:
 - Dữ liệu sau khi phân chia càng đồng nhất càng tốt
 - Đo bằng độ tăng thông tin (Information Gain - IG)
 - **Chọn thuộc tính có độ tăng thông tin lớn nhất**
 - IG dựa trên entropy của tập (con) dữ liệu

Entropy

- Entropy là một đại lượng được sử dụng trong lĩnh vực Lý thuyết thông tin
- Được sử dụng để đánh giá mức độ hỗn tạp của một tập dữ liệu
- Trường hợp dữ liệu S có 2 loại nhãn:

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

Trong đó: p_1 : tỉ lệ các mẫu trong tập S thuộc vào lớp 1, p_2 : tỉ lệ các mẫu trong tập S thuộc vào lớp 2

Entropy

- Trường hợp tổng quát: Khi tập dữ liệu S có C loại nhãn (C phân lớp):

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

Trong đó p_i là tỉ lệ các ví dụ trong tập S thuộc vào lớp i

Entropy - Ví dụ

- S gồm 14 ví dụ, trong đó có 9 ví dụ thuộc lớp c_1 và 5 ví dụ thuộc lớp c_2
- Entropy của tập S với 2 lớp:

$$Entropy(S) = -\frac{9}{14} \cdot \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \cdot \log_2\left(\frac{5}{14}\right) = 0.94$$

- Entropy=0, nếu tất cả các ví dụ cùng một lớp (c_1 hoặc c_2)
- Entropy=1, số lượng ví dụ thuộc về lớp c_1 bằng số lượng ví dụ thuộc về lớp c_2
- Entropy thuộc khoảng (0,1) nếu như số lượng các ví dụ thuộc về lớp c_1 khác với ví dụ thuộc về lớp c_2

Độ tăng thông tin - Information Gain

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

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

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

- Trong công thức trên, thành phần thứ 2 thể hiện giá trị Entropy sau khi tập S được phân chia bởi các giá trị thuộc tính A

Độ tăng thông tin - Ví dụ

- Hãy tính giá trị độ tăng thông tin của thuộc tính **Gió** đối với tập dữ liệu S

| Ngày | Trời | Nhiệt độ | Độ ẩm | Gió | Chơi tennis |
|------|------|------------|-------------|------|-------------|
| D1 | nắng | nóng | cao | yếu | không |
| D2 | nắng | nóng | cao | mạnh | không |
| D3 | u ám | nóng | cao | yếu | có |
| D4 | mưa | trung bình | cao | yếu | có |
| D5 | mưa | lạnh | bình thường | yếu | có |
| D6 | mưa | lạnh | bình thường | mạnh | không |
| D7 | u ám | lạnh | bình thường | mạnh | có |
| D8 | nắng | trung bình | cao | yếu | không |
| D9 | nắng | lạnh | bình thường | yếu | có |
| D10 | mưa | trung bình | bình thường | yếu | có |
| D11 | nắng | trung bình | bình thường | mạnh | có |
| D12 | u ám | trung bình | cao | mạnh | có |
| D13 | u ám | nóng | bình thường | yếu | có |
| D14 | mưa | trung bình | cao | mạnh | không |

Độ tăng thông tin - Ví dụ

Hãy tính giá trị độ tăng thông tin của thuộc tính **Gió** đối với tập dữ liệu S

- Tập dữ liệu S có 2 phân lớp: có (+), không (-)
- Thuộc tính Gió có 2 giá trị: yếu, mạnh
- $S = [9+, 5-]$
- $S_{\text{yếu}} = [6+, 2-]$
- $S_{\text{mạnh}} = [3+, 3-]$

$$\begin{aligned}
 IG(S, \text{Gió}) &= Entropy(S) - \sum_{v \in \{\text{yếu}, \text{mạnh}\}} \frac{|S_v|}{|S|} Entropy(S_v) \\
 &= Entropy(S) - \frac{8}{14} \cdot Entropy(S_{\text{yếu}}) - \frac{6}{14} \cdot Entropy(S_{\text{mạnh}}) \\
 &= 0.048
 \end{aligned}$$

Hoc cây quyết định - Ví dụ

- Tại nút gốc, thuộc tính nào trong số các thuộc tính {Trời, Nhiệt độ, Độ ẩm, Gió} nên được chọn là thuộc tính kiểm tra?
- Tính:
 - $IG(S, \text{Trời}) = 0.248$
 - $IG(S, \text{Nhiệt độ}) = 0.029$
 - $IG(S, \text{Độ ẩm}) = 0.151$
 - $IG(S, \text{Gió}) = 0.048$
- Thuộc tính **Trời** có giá trị IG cao nhất. Vì vậy, thuộc tính Trời được chọn làm thuộc tính kiểm tra cho nút gốc!

Học cây quyết định - Ví dụ

- Tại Node1, thuộc tính nào trong số các thuộc tính {Nhiệt độ, Độ ẩm, Gió} nên được chọn làm thuộc tính kiểm tra tiếp theo?
- Lưu ý: Thuộc tính Trời bị loại ra vì nó đã được sử dụng bởi cha của nút Node1 (là nút gốc).
- Tính:
 - $IG(S_{\text{nắng}}, \text{Nhiệt độ}) = 0.57$
 - $IG(S_{\text{nắng}}, \text{Độ ẩm}) = 0.97$
 - $IG(S_{\text{nắng}}, \text{Gió}) = 0.019$
- Vì vậy, thuộc tính **Độ ẩm** được chọn là thuộc tính kiểm tra cho nút Node1

Bài tập 1

Cho dữ liệu huấn luyện như trong bảng. Màu, Loại, Hãng là các thuộc tính, f là nhãn phân loại.

| Màu | Loại | Hãng | f |
|-------|-------|--------|---|
| Trắng | 7 chỗ | Toyota | - |
| Đen | 7 chỗ | Honda | + |
| Trắng | 5 chỗ | Honda | - |
| Đen | 5 chỗ | Toyota | + |
| Đỏ | 7 chỗ | Honda | + |
| Đỏ | 5 chỗ | Honda | - |
| Trắng | 5 chỗ | Toyota | + |

Hãy xác định nút gốc cho cây quyết định sử dụng thuật toán ID3. Trong trường hợp có nhiều thuộc tính có cùng mức độ ưu tiên thì chọn theo thứ tự từ trái sang phải (Màu, Loại, Hãng).

Các đặc điểm của ID3

- ID3 là thuật toán tìm kiếm cây quyết định phù hợp với dữ liệu huấn luyện
- Tìm kiếm theo kiểu tham lam, bắt đầu từ cây rỗng
- Hàm đánh giá là độ tăng thông tin
- ID3 có khuynh hướng (bias) lựa chọn cây đơn giản
 - Ít nút
 - Các thuộc tính có độ tăng thông tin lớn nằm gần gốc

Training error và Test error (1/2)

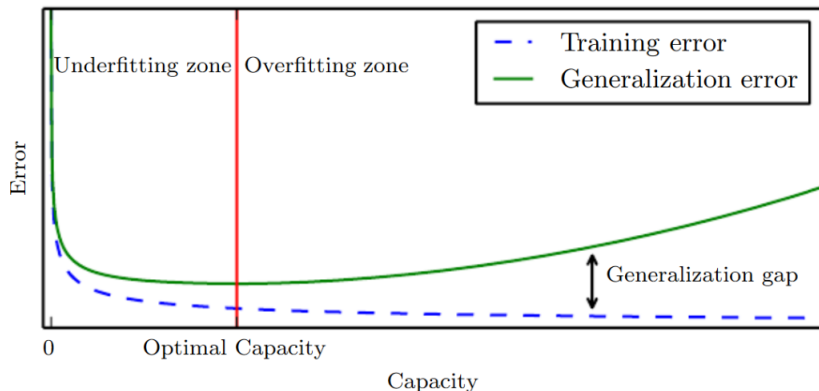
- Training error (lỗi huấn luyện)
 - Là lỗi đo được trên tập **dữ liệu huấn luyện**
 - Thường đo bằng sự sai khác giữa giá trị tính toán của mô hình và giá trị thực của dữ liệu huấn luyện
 - Trong quá trình học ta cố gắng làm **giảm tới mức tối thiểu lỗi huấn luyện**
- Test error (lỗi kiểm tra)
 - Là lỗi đo được trên tập **dữ liệu kiểm tra**
 - Là cái ta thực sự quan tâm

Làm sao ta có thể tác động tới hiệu quả của mô hình trên tập dữ liệu kiểm tra khi ta chỉ quan sát được tập dữ liệu huấn luyện?

Vấn đề quá vừa dữ liệu (overfitting)

- Quá vừa dữ liệu (data overfitting hay overfitting) là vấn đề thường gặp trong học máy và ảnh hưởng nhiều tới độ chính xác của các mô hình.
- Khi xây dựng cây quyết định, thuật toán cố gắng xây dựng cây phù hợp với dữ liệu một cách tối đa.
- Khi cây cho độ chính xác tốt trên dữ liệu huấn luyện nhưng lại cho kết quả không tốt trên dữ liệu kiểm tra => Cây quyết định quá vừa (overfitting) với dữ liệu huấn luyện.

Underfitting và Overfitting



Underfitting: dưới vừa; Overfitting: quá vừa

Generalization error = test error

Capacity: Khả năng của mô hình

Chống quá vừa bằng cách tỉa cây

- Chia dữ liệu thành hai phần
 - Huấn luyện
 - Kiểm tra
- Tạo cây đủ lớn trên dữ liệu huấn luyện
- Tính độ chính xác của cây trên tập kiểm tra
- Loại bỏ cây con sao cho kết quả trên dữ liệu kiểm tra được cải thiện nhất
- Lặp cho đến khi không còn cải thiện được kết quả nữa

Thuật toán C4.5

Thuật toán C4.5 là cải tiến của thuật toán ID3, thực hiện tĩa các luật như sau:

- Xây dựng cây quyết định cho phép phân loại đúng tối đa dữ liệu huấn luyện
- Biến đổi các cây thành các luật
- Tĩa từng luật bằng cách bỏ bớt các điều kiện thành phần nếu sau khi bỏ độ chính xác tăng lên.
- Sắp xếp các luật sau khi tĩa theo mức độ chính xác trên tập kiểm tra

Sử dụng các thuộc tính có giá trị liên tục

- Tạo ra những thuộc tính rời rạc mới
- Ví dụ, với những thuộc tính liên tục A , tạo ra những thuộc tính A_c như sau:
 - $A_c = \text{true}$ nếu $A > c$
 - $A_c = \text{false}$ nếu $A \leq c$
- Xác định ngưỡng c thế nào?: Thường chọn sao cho A_c đem lại độ tăng thông tin lớn nhất
- Có thể chia thành nhiều khoảng với nhiều ngưỡng

Ví dụ

- Chẳng hạn, nhiệt độ được cho dưới dạng số đo thực như trong ví dụ sau (ở đây nhiệt độ tính bằng độ F):

| | | | | | | |
|-------------|-------|-------|----|----|----|-------|
| Nhiệt độ | 45 | 56 | 60 | 74 | 80 | 90 |
| Chơi tennis | không | không | có | có | có | không |

- Xác định những trường hợp hai ví dụ nằm cạnh nhau nhưng có nhãn khác nhau.
- Giá trị trung bình của thuộc tính A của hai thuộc tính như vậy sẽ được sử dụng làm giá trị dự kiến của ngưỡng c

Ví dụ: $(56+60)/2 = 58$; $(80+90)/2 = 85$

- Tính độ tăng thông tin cho từng giá trị dự kiến và chọn c đem lại độ tăng thông tin lớn nhất (Nhiệt độ₅₈ và Nhiệt độ₈₅)

Phân loại Bayes đơn giản

Học dựa trên ví dụ