

Học máy

Giới thiệu

Lê Thành Sách

✉ Itsach@hcmut.edu.vn

Khoa Khoa học & Kỹ thuật Máy tính
Trường Đại học Bách Khoa - ĐHQG Tp.HCM

Tp.HCM. Ngày 1 tháng 9 năm 2019

Giới thiệu

Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

Mục lục

① Khái niệm

- Học máy
- Dữ liệu
- Dữ liệu - Học máy

② Các dạng học

③ Các cách tiếp cận khi phân tích dữ liệu

- Hai giai đoạn

④ Các công việc khi phân tích dữ liệu

- Chuẩn bị dữ liệu
- Xây dựng mô hình
- Lựa chọn mô hình và siêu tham số



Giới thiệu

1 Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu



Giới thiệu

Mục lục

2 **Khái niệm**

Học máy

Dữ liệu

Dữ liệu - Học máy

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

Khái niệm

Khái niệm

Học máy



Giới thiệu

Mục lục

Khái niệm

3

Học máy

Dữ liệu

Dữ liệu - Học máy

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

• Học máy:

- là lĩnh vực nghiên cứu giúp cho máy tính học¹ để đưa ra cách thực hiện một công việc nào đó thay cho việc yêu cầu máy tính giải quyết công việc qua dãy lệnh cụ thể².
- Ví dụ:
 - Dùng học máy để định danh người dùng
 - Dùng học máy để nhận dạng và tổng hợp tiếng nói
 - Dùng học máy để phân loại văn bản
 - Dùng học máy để dự báo nhu cầu về năng lượng, thực phẩm, v.v.

¹ học từ kinh nghiệm trong quá khứ \equiv học từ dữ liệu

² nhiều trường hợp, nếu không dùng học máy cũng không có giải thuật nào khác để thay thế!

• Dữ liệu:

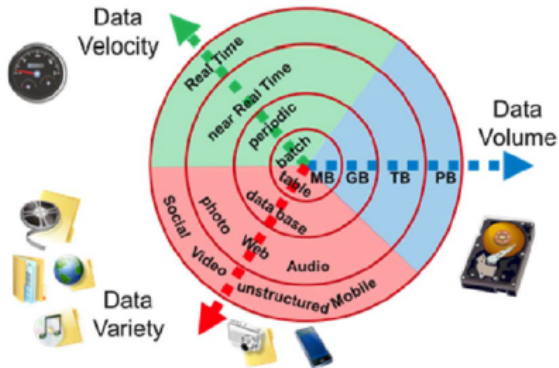
- là tập các giá trị biểu diễn cho các đo đạc, quan sát, và mô tả về một tín hiệu, quá trình hay hoạt động nào đó
- Ví dụ: video và hình ảnh trên Youtube, hồ sơ bệnh nhân trong bệnh viện, các tin nhắn người dùng, v.v.

• Dữ liệu lớn:

- là dữ liệu nhưng có các tính chất sau:
 - ❶ Độ lớn (**volume**): lượng dữ liệu là lớn tính theo dung lượng lưu trữ
 - ❷ Tốc độ (**velocity**): dữ liệu được sản sinh liên tục và cần được xử lý liên tục; ví dụ, dữ liệu của các camera giám sát
 - ❸ Tính đa dạng (**variety**): sự đa dạng về định dạng: văn bản, âm thanh, hình ảnh, chuỗi thời gian, v.v.

Khái niệm

Dữ liệu



Hình 1.1: Minh họa về dữ liệu lớn¹

¹nguồn: https://en.wikipedia.org/wiki/Big_data

Giới thiệu

Mục lục

Khái niệm

Học máy

5

Dữ liệu

Dữ liệu - Học máy

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

Giới thiệu

Dữ liệu - Học máy

Quan hệ hỗ tương

① Học máy cần dữ liệu để huấn luyện và chọn lựa mô hình, chọn lựa siêu tham số. Chọn kỹ thuật học máy nào là tùy vào mục đích ¹, sự sẵn sàng và tính chất của dữ liệu. Ví dụ:

- Dữ liệu lớn: nên chọn kỹ thuật xử lý từng bó nhỏ dữ liệu thay cho phải phối hợp tất cả các điểm dữ liệu;
- Dữ liệu ít có cấu trúc và quan hệ² đơn giản, không nên chọn mô hình quá phức tạp, như mạng nơron học sâu;
- Dữ liệu có yếu tố thời gian → chọn mô hình có yếu tố thời gian.

¹ mục đích: để phân loại, hồi quy; để gom nhóm trên đặc tính dữ liệu; hay để chọn lựa hành động phù hợp với ngữ cảnh

² quan hệ giữa đầu vào và đầu ra của bài toán



Giới thiệu

Mục lục

Khái niệm

Học máy

Dữ liệu

6

Dữ liệu - Học máy

Các dạng học

Các cách tiếp cận
khí phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

Giới thiệu

Dữ liệu - Học máy



Giới thiệu

Mục lục

Khái niệm

Học máy

Dữ liệu

7

Dữ liệu - Học máy

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

Quan hệ hỗ tương

② **Dữ liệu cần học máy để phân tích.** Thu thập dữ liệu gì và như thế nào là tùy vào yêu cầu của bài toán và kỹ thuật học máy được chọn. Ví dụ:

- Đã chọn kỹ thuật học có giám sát thì cần làm nhãn ¹ cho dữ liệu;
- Nếu kỹ thuật học giả thiết rằng các điểm dữ liệu được lấy mẫu độc lập thì việc chọn mẫu cũng tuân theo giả thiết đó;
- Yêu cầu của bài toán là định danh đến mức cá thể → phải gán nhãn là id của cá thể thay cho tên loại của cá thể.

¹từ gốc là **label**, hay **target**: là đầu ra của mô hình dự báo



Giới thiệu

Mục lục

Khái niệm

8 Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

Các dạng học

① Học có giám sát (**Supervised learning**)

- Đầu vào là dữ liệu \mathbf{X} và nhãn \mathbf{t}
 - \mathbf{X} : ma trận, kích thước: $N \times M$
 - \mathbf{t} : vectơ, kích thước N
 - mỗi hàng thứ n của \mathbf{t} là nhãn của điểm dữ liệu ở hàng n trong \mathbf{X}
- Mục đích: tìm ra một mô hình $f_w(\mathbf{x})^1$ từ tập học $\langle \mathbf{X}, \mathbf{t} \rangle$ để $f_w(\mathbf{x})$ dự báo được nhãn cho điểm dữ liệu mới \mathbf{x} là tốt nhất².
- Một số dạng bài toán quan trọng:
 - Hồi quy (**Regression**)
 - Phân loại (**Classification**)
 - Định danh, nhận dạng (**Identification, Recognition**)
 - Phát hiện (**Detection**)
 - Phân đoạn (**Segmentation**)
 - Hỏi-đáp (**Question-Answering**)

¹bài toán quy về làm tìm w ; w : tham số của mô hình; nếu w không có phần tử nào, $f_w(\mathbf{x})$ được gọi là mô hình phi tham số

²tính trên tập dữ liệu kiểm tra và sử dụng một độ đo đánh giá nào đó

② Học không giám sát (**Unsupervised learning**)

- Đầu vào là dữ liệu **X** , không có nhãn đi kèm
 - **X** : ma trận, kích thước: $N \times M$
- Mục đích: tìm ra những đặc tính quan trọng¹ trên dữ liệu đầu vào.
- Một số dạng bài toán quan trọng:
 - Gom nhóm (**Clustering**)
 - Thu giảm số chiều (**Dimensionality reduction**)
 - Trực quan hóa (**Visualization**)
 - Phát hiện bất thường (**Anomaly detection**)
 - Tìm kiếm luận kết hợp (**Association rule mining**)

¹tùy vào tiêu chí: sự tương quan giữa các điểm dữ liệu, tính khả tách của dữ liệu sau biến đổi, những điểm dữ liệu dị biệt so với nhóm, v.v.

③ Học tăng cường (**Reinforcement learning**)

- Đầu vào:
 - Môi trường của bài toán: không gian trên đó có thể thực hiện các hành động;
 - Danh sách các trạng thái của bài toán: ngữ cảnh cụ thể để thực hiện hành động;
 - Danh sách các hành động: một dạng tương tác với môi trường, có thể thay đổi tình trạng của môi trường;
 - Tiêu chí để phản hồi (dạng điểm) khi đối tượng thực hiện một hành động cụ thể nào đó.
- Mục đích: Thực hiện “thử-sai” nhiều lần để khi học xong thì có thể chọn lựa hành động tối ưu trên ngữ cảnh (trạng thái) cụ thể.
 - Trò chơi (**Games**)
 - Robotics (**Người-máy, xe và thiết bị vận hành/bay tự động**)
 - Thương mại (**Trading**)

Giới thiệu

Mục lục

Khái niệm

Các dạng học

12

Các cách tiếp cận
khí phân tích dữ
liệu

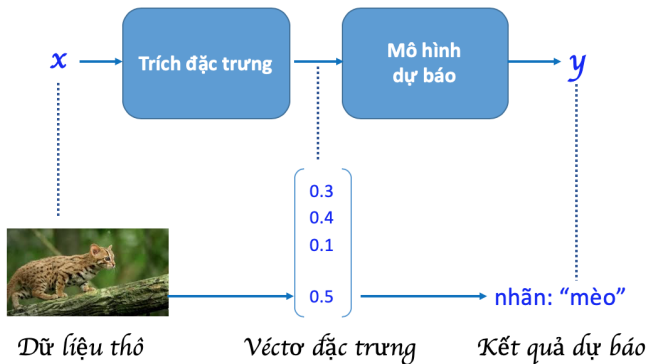
Hai giai đoạn

Các công việc khi
phân tích dữ liệu

Các cách tiếp cận khi phân tích dữ liệu

Các cách tiếp cận khi phân tích dữ liệu

Hai giai đoạn



Hình 3.1: Cách tiếp cận hai giai đoạn¹

¹Giai đoạn 1: rút trích đặc trưng từ dữ liệu thô; giai đoạn 2: sử dụng học máy để xây dựng mô hình dự báo

Giới thiệu

Mục lục

Khái niệm

Các dạng học

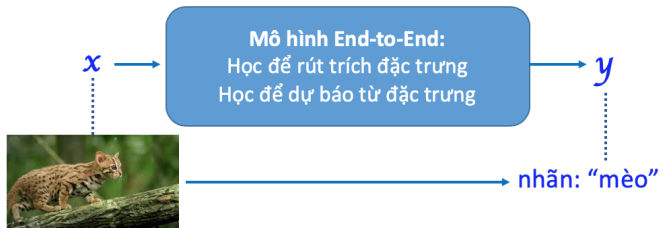
Các cách tiếp cận
khi phân tích dữ
liệu

13 Hai giai đoạn

Các công việc khi
phân tích dữ liệu

Các cách tiếp cận khi phân tích dữ liệu

Một giai đoạn



Hình 3.2: Cách tiếp cận một giai đoạn, **end-to-end** ¹

¹Đây là cách tiếp cận hiện đang phổ biến, mô hình học máy là mạng nơ-ron học sâu



Giới thiệu

Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

15

**Các công việc khi
phân tích dữ liệu**

Chuẩn bị dữ liệu

Xây dựng mô hình

Lựa chọn mô hình và siêu
tham số

24

ltsach@hcmut.edu.vn
Lê Thành Sách

Các công việc khi phân tích dữ liệu

Các công việc khi phân tích dữ liệu

- 1 Chuẩn bị dữ liệu
- 2 Xây dựng mô hình
- 3 Lựa chọn mô hình và siêu tham số
- 4 Triển khai ứng dụng



Giới thiệu

Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

16

**Các công việc khi
phân tích dữ liệu**

Chuẩn bị dữ liệu

Xây dựng mô hình

Lựa chọn mô hình và siêu
tham số

24

ltsach@hcmut.edu.vn
Lê Thành Sách

Các công việc khi phân tích dữ liệu

Chuẩn bị dữ liệu

Lưu ý

Để phân tích dữ liệu hiệu quả cần hiểu rõ về tính chất dữ liệu và có cảm nhận về dữ liệu!



Giới thiệu

Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

17 Chuẩn bị dữ liệu

Xây dựng mô hình

Lựa chọn mô hình và siêu
tham số

Các công việc khi phân tích dữ liệu

Chuẩn bị dữ liệu



Giới thiệu

Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

18 Chuẩn bị dữ liệu

Xây dựng mô hình

Lựa chọn mô hình và siêu
tham số

Các đầu việc

- Thu thập, làm nhẵn, làm sạch dữ liệu
- Thống kê và tóm tắt dữ liệu
- Rút trích đặc trưng
- Biến đổi và chuẩn hóa dữ liệu
- Hiển thị trực quan dữ liệu

Các công việc khi phân tích dữ liệu

Chuẩn bị dữ liệu



Giới thiệu

Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

19 Chuẩn bị dữ liệu

Xây dựng mô hình

Lựa chọn mô hình và siêu
tham số

Các đầu việc

- Chia dữ liệu thành các tập con:
 - Nếu dùng phương pháp đánh giá **hold-out**, 3 tập: Tập huấn luyện (training set), tập kiểm thử (validation set), và tập kiểm tra (test set)
 - Nếu dùng phương pháp đánh giá **cross-validation** hay **leave-one-out**, 2 tập: Tập huấn luyện và kiểm thử, tập kiểm tra

Các công việc khi phân tích dữ liệu

Chuẩn bị dữ liệu



Giới thiệu

Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận khi phân tích dữ liệu

Các công việc khi phân tích dữ liệu

20

Chuẩn bị dữ liệu

Xây dựng mô hình

Lựa chọn mô hình và siêu tham số

Lưu ý

- Không được đặt các điểm dữ liệu trong tập kiểm tra vào các tập kiểm thử và kiểm tra
- Cố tình thực hiện việc trên là một hình thức gian lận, nhằm có chỉ số đánh giá tốt cho mô hình

24

Các công việc khi phân tích dữ liệu

Xây dựng mô hình



Giới thiệu

Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

Chuẩn bị dữ liệu

21

Xây dựng mô hình

Lựa chọn mô hình và siêu
tham số

Các đầu việc

- Lựa chọn kỹ thuật học máy
 - Mô hình tuyến tính
 - Máy vectơ hỗ trợ
 - Mạng nơron/học sâu, v.v.
- Huấn luyện và kiểm thử mô hình để thu được mô hình tốt nhất
- Kiểm tra mô hình thu được và ghi nhận/công bố/tuyên bố kết quả

24

Các công việc khi phân tích dữ liệu

Lựa chọn mô hình và siêu tham số



Giới thiệu

Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

Chuẩn bị dữ liệu

Xây dựng mô hình

22

Lựa chọn mô hình và siêu
tham số

Khái niệm

- Mô hình được định nghĩa bởi:
 - Kỹ thuật học máy, ví dụ:
 - Mô hình tuyến tính
 - Máy vectơ hỗ trợ
 - Mạng nơron/học sâu, v.v.
 - Số lượng tham số, kiến trúc của các lớp tính toán (trong mạng nơron/học sâu)

24

Các công việc khi phân tích dữ liệu

Lựa chọn mô hình và siêu tham số



Giới thiệu

Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

Chuẩn bị dữ liệu

Xây dựng mô hình

23

Lựa chọn mô hình và siêu
tham số

Khái niệm

- Siêu tham số (**hyper-parameter**): là tham số chọn trước khi huấn luyện mô hình, quá trình huấn luyện không cho ra giá trị của nó: ví dụ:
 - Số tham số của mô hình hồi quy tuyến tính;
 - Hệ số **weight decay**
 - Kiến trúc và số lượng tham số trong mỗi lớp tính toán của mạng nơron

24

Các công việc khi phân tích dữ liệu

Lựa chọn mô hình và siêu tham số



Các đầu việc

- Chọn ra danh sách các tổ hợp của các siêu tham số
- Huấn luyện mô hình với từng bộ siêu tham số, trên tập huấn luyện
- Đánh giá mô hình thu được trên tập kiểm thử để cho ra các giá trị đánh giá
- Chọn lựa một tổ hợp làm cho giá trị đánh giá là tốt nhất
 \Rightarrow mô hình hay siêu tham số tối ưu

Giới thiệu

Mục lục

Khái niệm

Các dạng học

Các cách tiếp cận
khi phân tích dữ
liệu

Các công việc khi
phân tích dữ liệu

Chuẩn bị dữ liệu

Xây dựng mô hình

24

Lựa chọn mô hình và siêu
tham số

24