

Nội dung môn học

- Lecture 1: Giới thiệu về Học máy và khai phá dữ liệu
- **Lecture 2: Thu thập và tiền xử lý dữ liệu**
- Lecture 3: Hồi quy tuyến tính (Linear regression)
- Lecture 4+5: Phân cụm
- Lecture 6: Phân loại và Đánh giá hiệu năng
- Lecture 7: dựa trên láng giềng gần nhất (KNN)
- Lecture 8: Cây quyết định và Rừng ngẫu nhiên
- Lecture 9: Học dựa trên xác suất
- Lecture 10: Mạng nơron (Neural networks)
- Lecture 11: Máy vector hỗ trợ (SVM)
- Lecture 12: Khai phá tập mục thường xuyên và các luật kết hợp
- Lecture 13: Thảo luận ứng dụng trong thực tế

ĐẶT VẤN ĐỀ

- Khai phá dữ liệu là một quá trình phân tích dữ liệu theo nhiều khía cạnh và tổng hợp nó lại để có được thông tin hữu ích hay tri thức.

ĐẶT VẤN ĐỀ

- Các bước của quá trình phát hiện tri thức gồm
 1. Thu thập, lựa chọn dữ liệu
 2. Tiền xử lý dữ liệu
 3. Chuyển đổi
 4. Khai phá dữ liệu
 5. Giải thích/Đánh giá

ĐẶT VẤN ĐỀ

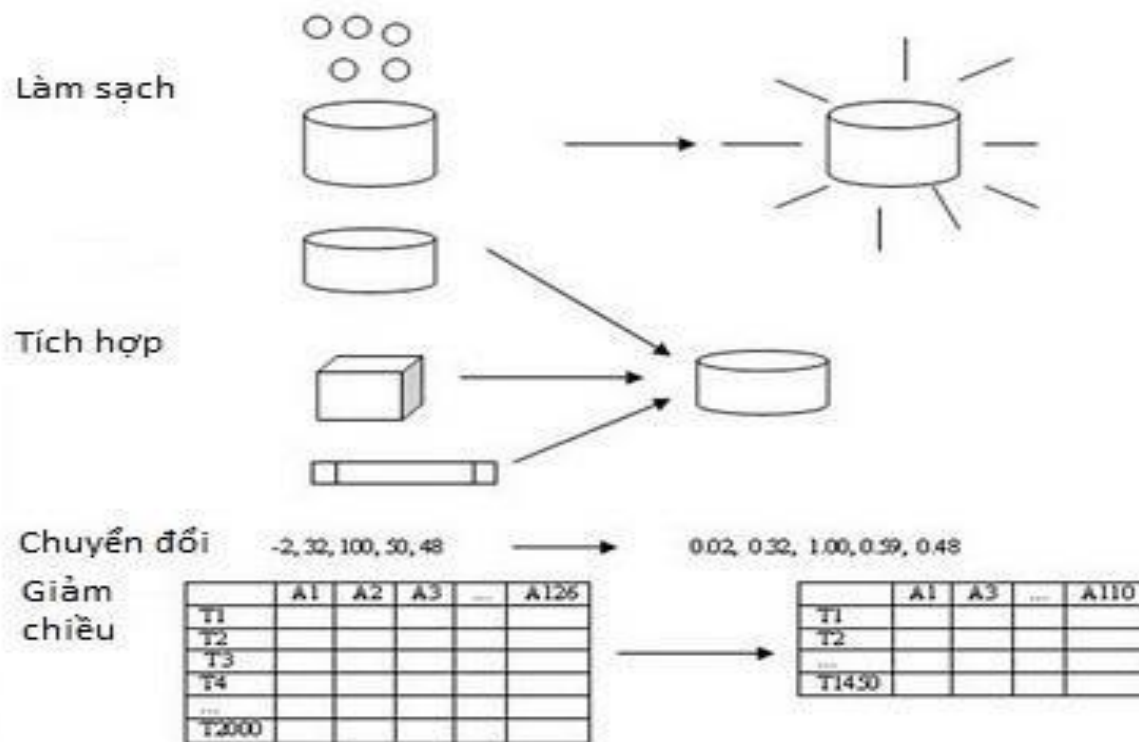
- Vì sao phải tiền xử lý dữ liệu ?
 - Không đầy đủ (Incomplete) : thiếu một vài giá trị thuộc tính
 - Nhiều (Noisy) : xuất hiện giá trị lỗi, lỗi chủ quan người nhập dữ liệu
 - Không nhất quán (Inconsistent) : sự khác biệt trong cách phân loại, phân biệt hay đơn vị của dữ liệu ...

ĐẶT VẤN ĐỀ

- Quy trình tiền xử lý dữ liệu
 1. **Làm sạch** : Loại bỏ các giá trị sai, kiểm tra tính nhất quan của dữ liệu.
 2. **Tích hợp** : Dữ liệu có nhiều nguồn nên cần lưu theo một cách thức thống nhất.
 3. **Chuyển đổi** : Chuẩn hóa và tập hợp dữ liệu.
 4. **Giảm chiều** : Mô tả dữ liệu trong kích thước nhỏ nhưng không làm mất kết quả cần kết xuất.

ĐẶT VẤN ĐỀ

- Quy trình tiền xử lý dữ liệu



MỤC LỤC

- Đặt vấn đề
- Làm sạch dữ liệu
- Tích hợp
- Giảm chiều

LÀM SẠCH DỮ LIỆU

- Đây là thủ tục quan trọng gồm ba bước chính
 1. Điền đầy các giá trị bị mất
 2. Chuốt dữ liệu để loại nhiễu
 3. Kiểm tra và sửa tính không nhất quán



LÀM SẠCH DỮ LIỆU

- **Bước 1:** Điền đầy các giá trị bị mất, có thể chọn một trong các phương pháp
 - Bỏ không xét đến bộ dữ liệu bị mất giá trị
 - Điền lại giá trị bằng tay
 - Gán cho giá trị nhãn đặc biệt hay ngoài khoảng biểu diễn
 - Gán giá trị trung bình cho nó.
 - Gán giá trị trung bình của các mẫu khác thuộc cùng lớp đó.
 - Tìm giá trị có xác suất lớn nhất điền vào chỗ bị mất (**hồi quy, suy diễn Bayes, cây quyết định qui nạp**)

LÀM SẠCH DỮ LIỆU

- **Bước 2:** Chuốt dữ liệu loại nhiễu, có thể chọn một trong các phương pháp
 - Hồi quy (Regression) : sẽ dành chương riêng
 - Phân cụm (Cluster) : sẽ dành chương riêng

LÀM SẠCH DỮ LIỆU

- **Bước 3: kiểm tra và sửa tính** không nhất quán trong dữ liệu.
 - Đề phát hiện kiểm tra sự bất thường trong giá trị dữ liệu, ta tuân thủ ba luật sau
 - Luật giá trị duy nhất (unique rule) : mỗi giá trị của một thuộc tính sẽ khác biệt với các giá trị khác thuộc cùng thuộc tính.
 - Luật liên tục (consecutive rule) : không có giá trị bị mất giữa giá trị lớn nhất và nhỏ nhất tương ứng một thuộc tính
 - Luật giá trị rỗng (null rule) : xác định trước các ký hiệu hay cách đánh dấu giá trị rỗng
 - Dùng để sửa tính không nhất quán dữ liệu
 - Công cụ chà dữ liệu (Data scrubbing tools) : dùng cho một lĩnh vực cụ thể.
 - Công cụ kiểm toán dữ liệu (Data auditing tools) : dùng cho việc phân tích dữ liệu, xác định quan hệ, xác định các luật.

MỤC LỤC

- Đặt vấn đề
- Làm sạch dữ liệu
- Tích hợp
- Giảm chiều

TÍCH HỢP

- **Ý nghĩa tích hợp và chuyển đổi dữ liệu, chuẩn hóa để tiến hành khai phá/học máy**
 - Tập hợp : các giá trị dữ liệu tạo thành bộ hay khối
 - Tổng quát hóa (generalization) : các dữ liệu "thô" được thay bằng các khái niệm đã chuẩn hóa
 - Chuẩn hóa (normalization) : nếu phạm vi dữ liệu lớn thì đưa nó về phạm vi chuẩn
 - Xây dựng thuộc tính (attribute construction) : thuộc tính mới thêm vào giúp quá trình khai phá dữ liệu

MỤC LỤC

- Đặt vấn đề
- Làm sạch dữ liệu
- Tích hợp
- Giảm chiều

GIẢM CHIỀU

- **Ý nghĩa:** Việc giảm kích thước của dữ liệu cần đồng thời giữ được tính phân tích dữ liệu, tăng tốc quá trình khai phá/học máy

GIẢM CHIỀU

- Các chiến lược **giảm kích thước dữ liệu**
 - **Lựa chọn tập con các thuộc tính** : trong đó các thuộc tính không liên quan, dư thừa hoặc các chiều cũng có thể xóa hay loại bỏ
 - **Giảm chiều** : trong đó cơ chế mã hóa được sử dụng để giảm kích cỡ tập dữ liệu
 - **Rời rạc hóa và trừu tượng khái niệm** : trong đó các giá trị dữ liệu thô được thay thế bằng các khái niệm trừu tượng đã rời rạc hóa.

TỔNG KẾT

- Vấn đề khi tiến hành thu thập dữ liệu dùng cho bài toán khai phá dữ liệu/học máy.
- Cần theo các bước của quy trình thu thập và tiền xử lý dữ liệu.
- Cần hiểu ý nghĩa trong từng bước của quy trình.