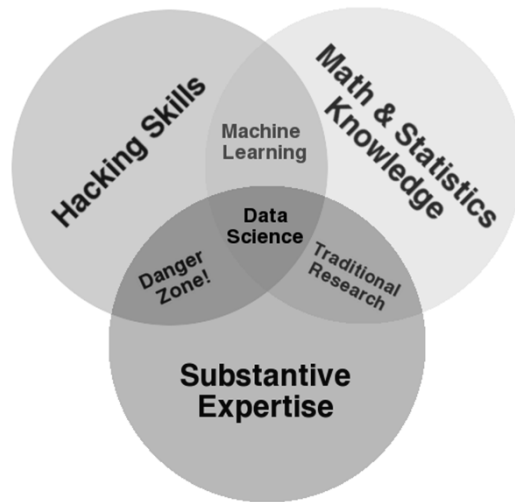


THƯ VIỆN TRONG PYTHON HỖ TRỢ PHÂN TÍCH DỮ LIỆU

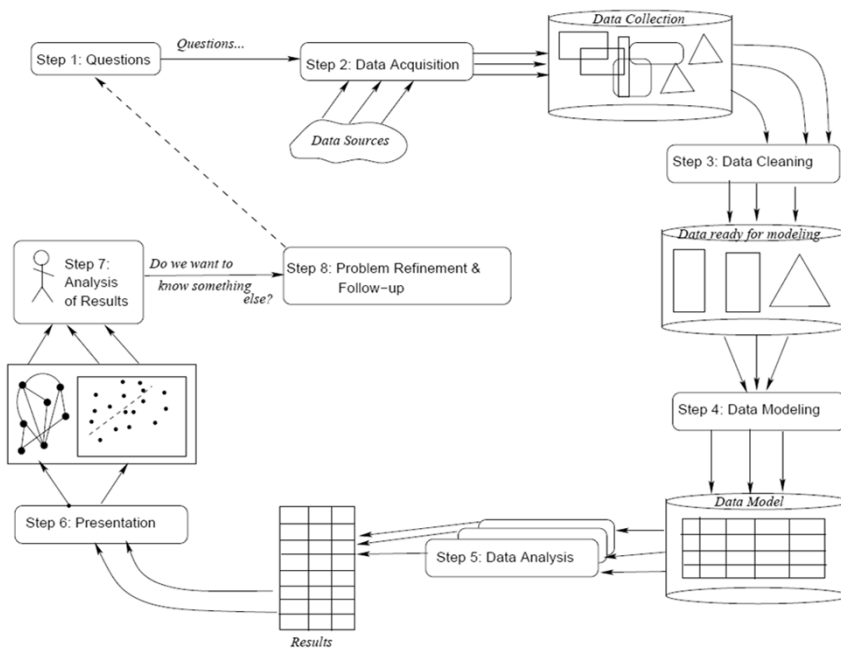
Nội dung

1. Thư viện scikit-learn

Khoa học Dữ liệu và máy Học



Quá trình xử lý của khoa học dữ liệu



Ví dụ: hệ thống phát hiện thư rác

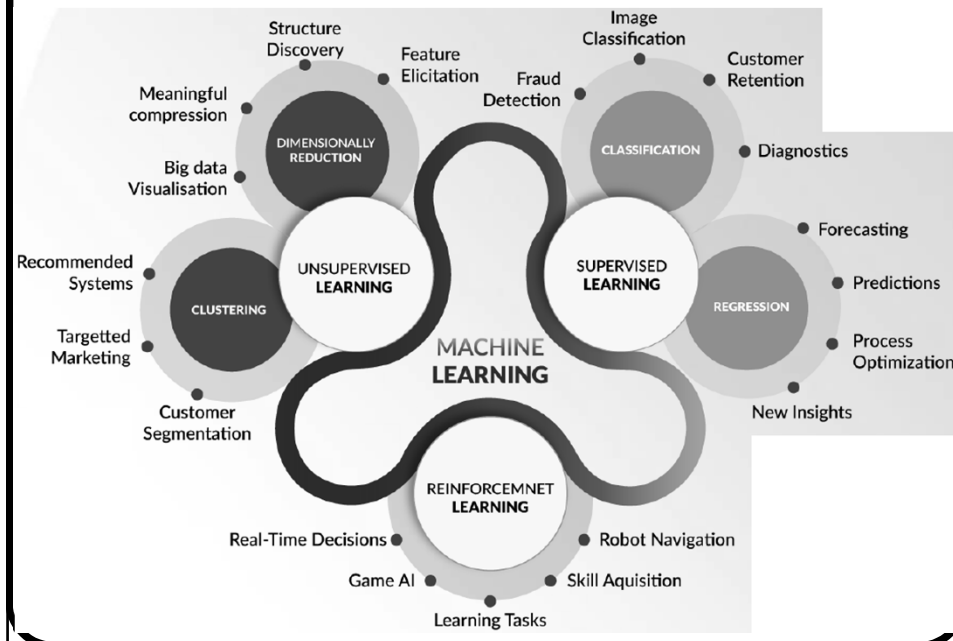
1. Thu thập mẫu thư (gồm cả thư rác và thư thường)
2. Xác định đề bài (phân lớp hay đánh giá)
3. Xử lý dữ liệu
4. Chọn mô hình học máy phù hợp với bài toán phân loại thư rác
5. Huấn luyện mô hình
6. Hiệu chỉnh, tinh chỉnh mô hình
7. Áp dụng thực tế (chạy trên email server thực)
8. Tiếp tục cập nhật theo phản hồi của người dùng

Học có giám sát (supervised learning)

Học không giám sát (unsupervised learning)

Học bán giám sát

Bài toán cơ bản

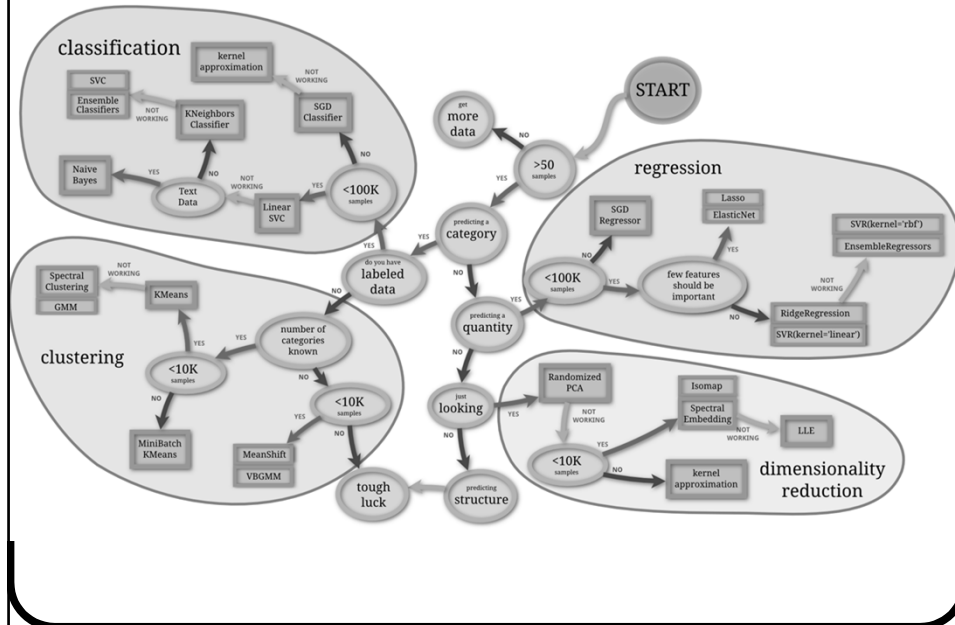


Thư viện học máy scikit-learn

Scikit-learn cung cấp gần như tất cả các loại thuật toán học máy cơ bản cùng với các kỹ thuật xử lý dữ liệu đã được chuẩn hóa

Cài đặt: **pip install scikit-learn scipy** hoặc
py- m pip install scikit-learn scipy

Chọn thuật toán học máy phù hợp



Các dạng học

1 Học có giám sát (Supervised learning)

- Đầu vào là dữ liệu X và nhãn t
 - X : ma trận, kích thước: $N \times M$
 - t : vectơ, kích thước N
 - mỗi hàng thứ n của t là nhãn của điểm dữ liệu ở hàng n trong X
- Mục đích: tìm ra một mô hình $f_w(x)$ ¹ từ tập học $\langle X, t \rangle$ để $f_w(x)$ dự báo được nhãn cho điểm dữ liệu mới 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, $f_w(x)$ 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**)

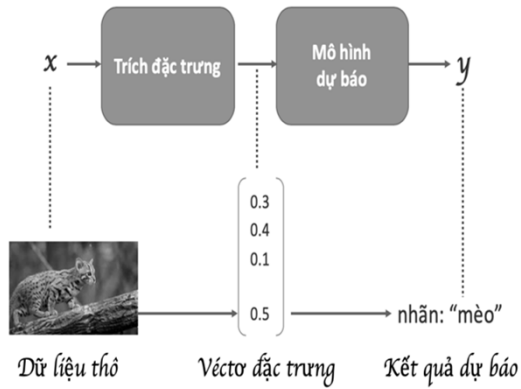
¹tùy vào tiêu chí: sự tương quan giữa các điểm dữ liệu, tính khả tách của dữ liệu sau biến đổi, những điểm dữ liệu dị biệt so với nhóm, v.v.

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**)

Các cách tiếp cận khi phân tích dữ liệu

Hai giai đoạn

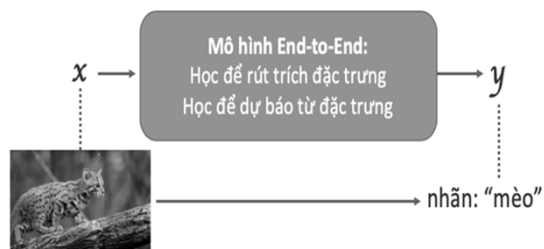


Hình 3.1: Cách tiếp cận hai giai đoạn¹

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

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

Các công việc khi phân tích dữ liệu

- ❶ Chuẩn bị dữ liệu
- ❷ Xây dựng mô hình
- ❸ Lựa chọn mô hình và siêu tham số
- ❹ Triển khai ứng dụng