

Biên soạn: TS. ĐỖ THANH NGHỊ - TS. PHẠM NGUYỄN KHANG

---

GIÁO TRÌNH

# NGUYÊN LÝ MÁY HỌC



**NHÀ XUẤT BẢN ĐẠI HỌC CẦN THƠ**  
2012

**BIÊN MỤC TRƯỚC XUẤT BẢN THỰC HIỆN BỞI  
TRUNG TÂM HỌC LIỆU TRƯỜNG ĐẠI HỌC CẦN THƠ**

Đỗ, Thanh Nghị  
Giáo trình nguyên lý máy học / Đỗ Thanh Nghị, Phạm Nguyên Khang .- Cần Thơ :  
Nxb. Đại học Cần Thơ, 2012  
154 tr. : minh họa ; 24 cm  
Sách có danh mục tài liệu tham khảo

1. Machine learning                      2. Image processing  
3. Nguyên lý máy học                    4. Giải thuật máy học

I. Nhan đề                                      II. Phạm, Nguyên Khang

006.31 - DDC 22  
Ngh300

MFN 175317

## LỜI GIỚI THIỆU

Nhằm góp phần làm phong phú thêm nguồn tư liệu phục vụ nghiên cứu, học tập cho bạn đọc trong và ngoài ngành Công nghệ thông tin, Nhà Xuất bản Đại học Cần Thơ xin được phép ấn hành và giới thiệu cùng bạn đọc quyển sách giáo trình “Nguyên lý máy học” do TS. Đỗ Thanh Nghị và TS. Phạm Nguyên Khang biên soạn. Giáo trình bao gồm 6 chương với 154 trang; được tổ chức thành 4 phần một cách có hệ thống. Phần 1 giới thiệu khái quát về máy học và nguyên lý học thống kê. Phần 2 trình bày các giải thuật học có giám sát tiêu biểu. Phần 3 trình bày các giải thuật học không giám sát. Phần 4 giới thiệu các ứng dụng phổ biến của máy học. Nội dung các chương giới thiệu nguyên lý học thống kê, mạng nơ-ron nhân tạo, máy học véc-tơ hỗ trợ (SVM), phương pháp đánh giá mô hình phân lớp, các giải thuật gom nhóm dữ liệu như *k-means*, bản đồ tự tổ chức (SOM), cực đại hóa kỳ vọng (EM) và ứng dụng nhận dạng ký tự số viết tay, phân lớp tự động văn bản. Thêm vào đó, cuối mỗi chương còn có rất nhiều bài tập và tài liệu tham khảo hữu ích cho bạn đọc. Giáo trình là tài liệu tham khảo có giá trị cho sinh viên ngành Công nghệ thông tin và bạn đọc muốn tìm hiểu và nghiên cứu về máy học.

Nhà Xuất bản Đại học Cần Thơ chân thành cảm ơn các Tác giả và sự đóng góp ý kiến của quý Thầy Cô trong Hội đồng thẩm định trường Đại học Cần Thơ để giáo trình “Nguyên lý máy học” được ra mắt bạn đọc.

Nhà Xuất bản Đại học Cần Thơ trân trọng giới thiệu đến giảng viên, sinh viên và bạn đọc giáo trình này.

Chân thành cảm ơn!

**NHÀ XUẤT BẢN ĐẠI HỌC CẦN THƠ**

## LỜI NÓI ĐẦU

Đã trải qua nhiều thập kỷ, chúng ta chứng kiến sự phát triển của ngành khoa học máy tính. Xuất phát điểm từ công nghệ chế tạo phần cứng, lập trình cấp thấp, đến công nghệ thông tin, tổ chức, quản lý và xử lý hiệu quả hệ thống thông tin. Hiện nay, ngành khoa học máy tính đang bước vào thời kỳ công nghệ tri thức, ở mức trù tượng cao hơn, làm máy tính trở nên thông minh hơn, để có thể giúp con người giải quyết nhiều vấn đề phức tạp trong thực tiễn. Chẳng hạn như máy học là hướng tiếp cận nhằm phát triển các kỹ thuật cho phép các máy tính có thể học từ dữ liệu để nhận dạng các mẫu phức tạp và đưa ra các quyết định thông minh. Máy học có tính ứng dụng rất cao trong thực tế như chẩn đoán y khoa, phát hiện thẻ tín dụng giả, phân tích thị trường chứng khoán, phân loại các chuỗi DNA, nhận dạng tiếng nói và chữ viết, dịch tự động, chơi trò chơi và điều khiển robot. Chính vì vậy, nguyên lý máy học trở thành môn học chuyên ngành rất quan trọng của chương trình đào tạo cử nhân, kỹ sư và thạc sĩ công nghệ thông tin. Trong thời gian giảng dạy vừa qua, chúng tôi nhận thấy có quá ít tài liệu tham khảo trình bày một cách có hệ thống về nguyên lý máy học. Chính vì lý do đó, chúng tôi đã thực hiện biên soạn quyển giáo trình “Nguyên lý máy học” nhằm mục đích cung cấp thêm tài liệu tham khảo bằng tiếng Việt đến độc giả là sinh viên chuyên ngành công nghệ thông tin.

Giáo trình được soạn dựa trên kiến thức, kinh nghiệm học tập nghiên cứu, giảng dạy của chúng tôi trong suốt từ năm 2000 đến nay và các tài liệu tham khảo chính được trình bày trong mục tài liệu tham khảo của từng chương. Chúng tôi cung cấp các kiến thức cơ bản và các giải thích nhằm đơn giản việc trình bày các giải thuật học. Cuối mỗi chương có bài tập giúp độc giả ôn tập lại kiến thức của chương và trao dồi kỹ năng thực hành. Giáo trình được tổ chức thành sáu chương:

Chương 1 tập trung vào việc giới thiệu các kiến thức cơ bản của máy học. Nội dung bao gồm khái niệm học có giám sát, học không giám sát, nguyên lý học thống kê được xem là nền tảng của các kỹ thuật học từ dữ liệu.

Chương 2 giới thiệu kiến thức cơ bản về giải thuật học có giám sát mạng nơ-ron nhân tạo. Chúng tôi trình bày về mô hình mạng nơ-ron của McCulloch và Pitts, kiến trúc mạng nơ-ron, giải thuật huấn luyện mạng Perceptron, mạng nơ-ron đa tầng cho vấn đề phân lớp.

Chương 3 giới thiệu về giải thuật học có giám sát máy học véc-tơ hỗ trợ cho vấn đề phân lớp, hồi quy và phát hiện phần tử cá biệt. Nội dung bắt đầu từ việc mô hình hóa bài toán đến giải thuật máy học véc-tơ hỗ trợ. Chương được kết thúc với phần thảo luận về giải thuật máy học véc-tơ hỗ trợ.

Chương 4 cung cấp cho sinh viên về phương pháp đánh giá hiệu quả của giải thuật học có giám sát. Các nghi thức kiểm tra và tiêu chí để so sánh hiệu quả của giải thuật học.

Chương 5 tập trung trình bày các giải thuật học không giám sát như giải thuật gom nhóm k-means, bản đồ tự tổ chức và giải thuật cực đại hóa kỳ vọng. Chúng tôi giới thiệu kiến thức cơ bản từ việc mô hình hóa bài toán đến giải thuật học cho gom nhóm.

Chương 6 giới thiệu hai ứng dụng phổ biến trong thực tế của máy học là nhận dạng chữ viết tay và phân loại văn bản tự động với mạng nơ-ron, máy học véc-tơ hỗ trợ.

Nội dung giáo trình được giảng dạy cho sinh viên kỹ sư, thạc sĩ công nghệ thông tin trong thời lượng 30 tiết lý thuyết. Giảng viên có thể phân bổ thời gian dạy lý thuyết kết hợp với các ứng dụng minh họa của máy học. Chương 1, 2, 3, 4 có thể được giảng dạy trong 18 tiết lý thuyết. Chương 5 có thể được trình bày trong 10 tiết lý thuyết, và thời gian còn lại là của chương 6. Nếu giảng dạy cho sinh viên kỹ sư, cần giảm bớt hai giải thuật bản đồ tự tổ chức và cực đại kỳ vọng cho gom nhóm.

Trong thời gian biên soạn tài liệu, chúng tôi nhận được sự hỗ trợ quý báu về vật chất, tinh thần từ Khoa Công Nghệ Thông Tin và Truyền Thông, Trường Đại Học Cần Thơ. Chúng tôi đặc biệt gửi lời cảm ơn đến Quý Thầy Cô thuộc Bộ Môn Khoa Học Máy Tính, Hội Đồng Thẩm Định Giáo Trình, đã nhiệt tình góp ý cho bản thảo.

Do thời gian có hạn và lần đầu tiên biên soạn, tài liệu không thể tránh khỏi những thiếu sót. Chúng tôi mong nhận được sự góp ý chân thành từ quý độc giả để quyển sách ngày được hoàn thiện hơn.

Cần Thơ, tháng 6 năm 2012

Đỗ Thanh Nghị - Phạm Nguyên Khang

## MỤC LỤC

<b>PHẦN I. GIỚI THIỆU</b>	<b>1</b>
<b>Chương 1. GIỚI THIỆU</b>	<b>3</b>
1.1 KHÁI NIỆM CƠ BẢN MÁY HỌC	3
1.2 NGUYÊN LÝ HỌC THỐNG KÊ	6
1.2.1 Mô hình ước lượng hàm	7
1.2.2 Cực tiểu hóa rủi ro (minimizing a risk functional)	7
1.2.3 Cực tiểu rủi ro thực nghiệm (empirical risk minimization)	8
1.2.4 Nguyên lý suy diễn thực nghiệm (empirical inference)	9
1.2.5 Nguyên lý quy nạp cực tiểu rủi ro cấu trúc (structural risk minimization induction principle)	11
1.2.6 Nguyên lý quy nạp cực tiểu rủi ro cấu trúc cho mạng nơ-ron và máy học véc-tơ hỗ trợ SVM	13
1.3 KẾT LUẬN	14
BÀI TẬP	15
TÀI LIỆU THAM KHẢO	15
<b>PHẦN II. GIẢI THUẬT HỌC CÓ GIÁM SÁT</b>	<b>17</b>
<b>Chương 2. MẠNG NƠ-RON NHÂN TẠO</b>	<b>19</b>
2.1 GIỚI THIỆU	19
2.2 MÔ HÌNH NƠ-RON CỦA MCCULLOCH VÀ PITTS	20
2.3 MÔ HÌNH PERCEPTRON ĐA TẦNG	24
2.3.1 Mô hình Perceptron	24
2.3.2 Huấn luyện Perceptron	27
2.3.3 Perceptron đa tầng	38
2.3.4 Huấn luyện mạng nơ-ron MLP	40
2.4 MỘT SỐ VẤN ĐỀ CẦN CHÚ Ý KHI SỬ DỤNG GIẢI THUẬT HUẤN LUYỆN MẠNG NƠ-RON MLP	45
2.4.1 Điều kiện dừng	45
2.4.2 Khởi tạo trọng số	45
2.4.3 Tốc độ học	45
2.4.4 Thêm quán tính (inertia) vào lượng cập nhật của các trọng số	45
2.4.5 Hàm kích hoạt	46

2.4.6	Tiền xử lý đầu ra	46
2.4.7	Tiền xử lý đầu vào	46
2.4.8	Dữ liệu nhiều hơn hai lớp	46
2.4.9	Sử dụng mạng nơ-ron giải bài toán hồi quy	46
2.5	KẾT LUẬN	47
	BÀI TẬP	47
	TÀI LIỆU THAM KHẢO	49
<b>Chương 3.</b>	<b>MÁY HỌC VÉC-TƠ HỖ TRỢ</b>	<b>51</b>
3.1	MÁY HỌC SVM CHO PHÂN LỚP	52
3.2	GIẢI THUẬT MÁY HỌC SVM	61
3.3	MÁY HỌC SVM CHO PHÂN LỚP ĐA LỚP	63
3.4	MÁY HỌC SVM CHO HỒI QUY	65
3.5	MÁY HỌC SVM CHO PHÁT HIỆN PHẦN TỬ CÁ BIỆT	69
3.6	KẾT LUẬN	71
	BÀI TẬP	71
	TÀI LIỆU THAM KHẢO	73
<b>Chương 4.</b>	<b>PHƯƠNG PHÁP ĐÁNH GIÁ HIỆU QUẢ GIẢI THUẬT HỌC CÓ GIÁM SÁT</b>	<b>75</b>
4.1	NGHI THỨC KIỂM TRA	75
4.2	ĐỘ ĐO HIỆU QUẢ CỦA GIẢI THUẬT	76
	TÀI LIỆU THAM KHẢO	80
<b>PHẦN III.</b>	<b>GIẢI THUẬT HỌC KHÔNG GIÁM SÁT</b>	<b>81</b>
<b>Chương 5.</b>	<b>HỌC KHÔNG GIÁM SÁT</b>	<b>83</b>
5.1	BÀI TOÁN GOM NHÓM DỮ LIỆU	83
5.2	GIẢI THUẬT k-MEANS	84
5.3	BẢN ĐỒ TỰ TỔ CHỨC (SOM)	95
5.3.1	Mạng nơ-ron Kohonen	95
5.3.2	Giải thuật xây dựng bản đồ tự tổ chức	96
5.4	ƯỚC LƯỢNG MẬT ĐỘ XÁC SUẤT VÀ CỰC ĐẠI HÓA KỶ VỌNG	101
5.4.1	Bài toán ước lượng mật độ xác suất	101
5.4.2	Ước lượng mật độ bằng phương pháp hợp lý cực đại	101
5.4.3	Giải thuật cực đại hóa kỳ vọng	105

5.5 GOM NHÓM DỮ LIỆU VỚI CỰC ĐẠI HÓA KỶ VỌNG	109
5.5.1 Mô hình hóa bài toán	109
5.5.2 Mô hình hỗn hợp các phân phối chuẩn	111
5.5.3 Quan hệ giữa k-means và mô hình hỗn hợp các phân phối chuẩn	115
5.5 KẾT LUẬN	115
BÀI TẬP	118
TÀI LIỆU THAM KHẢO	119
<b>PHẦN IV. ỨNG DỤNG</b>	<b>121</b>
<b>Chương 6. ỨNG DỤNG</b>	<b>123</b>
6.1 NHẬN DẠNG CHỮ VIẾT TAY	123
6.2 PHÂN LỚP DỮ LIỆU VĂN BẢN	130
6.3 KẾT LUẬN	135
BÀI TẬP	136
TÀI LIỆU THAM KHẢO	136



## DANH MỤC HÌNH

<b>Hình 1.1:</b> Bài toán phân lớp.....	5
<b>Hình 1.2:</b> Bài toán hồi quy.....	5
<b>Hình 1.3:</b> Bài toán gom nhóm dữ liệu.....	6
<b>Hình 1.4:</b> Rủi ro và rủi ro thực nghiệm .....	9
<b>Hình 1.5:</b> Ví dụ về chiều VC, tập siêu phẳng trong không gian 2 chiều, khả tách.....	11
<b>Hình 1.6:</b> Ví dụ về chiều VC, tập siêu phẳng trong không gian 2 chiều, không khả tách....	11
<b>Hình 1.7:</b> Các mô hình học để phân lớp dữ liệu.....	13
<b>Hình 2.1:</b> Mô hình nơ-ron của McCulloch và Pitts .....	21
<b>Hình 2.2:</b> Hàm mạng bậc 2 .....	22
<b>Hình 2.3:</b> Mạng nơ-ron không chu trình (a) và mạng nơ-ron có chu trình (b) .....	24
<b>Hình 2.4:</b> Mô hình perceptron.....	25
<b>Hình 2.5:</b> Mô hình Perceptron không ngưỡng với đầu vào giả $x_0 = 1$ .....	26
<b>Hình 2.6:</b> Siêu phẳng hai chiều (đường thẳng) của một Perceptron phân lớp dữ liệu thành hai lớp: dương (+), âm (o).....	27
<b>Hình 2.7:</b> Perceptron mô phỏng hàm logic AND.....	30
<b>Hình 2.8:</b> Kết quả của Perceptron mô phỏng hàm logic AND .....	33
<b>Hình 2.9:</b> Mạng nơ-ron MLP 3 tầng.....	39
<b>Hình 2.10:</b> Minh họa ánh xạ phi tuyến của mạng nơ-ron MLP.....	40
<b>Hình 2.11:</b> Tính đạo hàm riêng của nơ-ron b ở tầng đầu ra .....	41
<b>Hình 2.12:</b> Tính đạo hàm riêng theo nơ-ron b ở tầng bất kỳ.....	43
<b>Hình 3.1:</b> Dạng mô hình máy học véc-tơ hỗ trợ .....	51
<b>Hình 3.2:</b> Vấn đề phân lớp tuyến tính .....	52
<b>Hình 3.3:</b> Phân lớp tuyến tính với SVM.....	53
<b>Hình 3.4:</b> SVM phân lớp tuyến tính dữ liệu không tách rời.....	54
<b>Hình 3.5:</b> Vấn đề phân lớp phi tuyến .....	56
<b>Hình 3.6:</b> Phân lớp tuyến tính trong không gian trung gian .....	56
<b>Hình 3.7:</b> Mô hình phân lớp SVM tuyến tính với $c = 10000$ .....	58
<b>Hình 3.8:</b> Mô hình phân lớp SVM tuyến tính với $c = 2$ .....	59
<b>Hình 3.9:</b> Mô hình phân lớp SVM tuyến tính ( $t = 0$ và $c = 10000$ ) .....	60
<b>Hình 3.10:</b> Mô hình phân lớp SVM phi tuyến sử dụng hàm nhân đa thức bậc 5 ( $t = 1$ , $d = 5$ và $c = 100000$ ) .....	60
<b>Hình 3.11:</b> Mô hình phân lớp SVM phi tuyến sử dụng hàm nhân RBF ( $t = 2$ , $\gamma = 100$ và $c = 100000$ ).....	61

<b>Hình 3.12:</b> SVM đa lớp với 1-tất cả (trái), 1-1 (phải) .....	63
<b>Hình 3.13:</b> SVM đa lớp với chiến lược tách 2 nhóm có lợi .....	64
<b>Hình 3.14:</b> LibSVM đa lớp với 1-1 .....	64
<b>Hình 3.15:</b> Xử lý bài toán hồi quy bằng SVM.....	65
<b>Hình 3.16:</b> Hồi quy SVM tuyến tính ( $t = 0$ , $\varepsilon = 0.05$ và $c = 1000$ ).....	67
<b>Hình 3.17:</b> Hồi quy SVM phi tuyến sử dụng hàm nhân đa thức bậc 5 ( $t = 1$ , $d = 5$ , $\varepsilon = 0.05$ và $c = 100000$ ) .....	68
<b>Hình 3.18:</b> Hồi quy SVM phi tuyến sử dụng hàm nhân RBF ( $t = 2$ , $\gamma = 10$ , $\varepsilon = 0.05$ và $c = 100000$ ) .....	68
<b>Hình 3.19:</b> SVM cho bài toán 1 lớp (phát hiện phần tử cá biệt).....	69
<b>Hình 3.20:</b> Phát hiện phần tử cá biệt với SVM phi tuyến sử dụng hàm nhân RBF ( $t = 2$ , $\gamma = 20$ , $v = 0.3$ và $c = 100000$ ) .....	70
<b>Hình 4.1:</b> Nghi thức kiểm tra chéo 10-fold .....	75
<b>Hình 5.1:</b> Gom nhóm dữ liệu.....	84
<b>Hình 5.2:</b> Ví dụ minh họa giải thuật k-means .....	86
<b>Hình 5.3:</b> Đồ thị phân bố của tập dữ liệu mô phỏng .....	88
<b>Hình 5.4:</b> Kết quả gom nhóm dữ liệu mô phỏng bằng giải thuật k-means.....	92
<b>Hình 5.5:</b> Ma trận scatterplot 2 chiều của dữ liệu Iris .....	93
<b>Hình 5.6:</b> Hiển thị kết quả gom dữ liệu Iris thành 3 nhóm của k-means.....	94
<b>Hình 5.7:</b> Mô hình một nơ-ron trong mạng nơ-ron Kohonen.....	95
<b>Hình 5.8:</b> Sơ đồ tổ chức của một mạng nơ-ron Kohonen .....	96
<b>Hình 5.9:</b> Huấn luyện bản đồ tự tổ chức.....	98
<b>Hình 5.10:</b> Hiển thị kết quả gom nhóm dữ liệu Iris bằng SOM trên lưới 5x5 .....	101
<b>Hình 5.11:</b> Hiển thị kết quả gom dữ liệu Iris thành 3 nhóm của EM .....	114
<b>Hình 6.1:</b> Máy học nhận dạng ký tự số viết tay .....	123
<b>Hình 6.2:</b> Mẫu ký tự số viết tay .....	124
<b>Hình 6.3:</b> Kiến trúc mạng tích chập cho nhận dạng ký tự số viết tay.....	125
<b>Hình 6.4:</b> Kiến trúc mạng tích chập với 'C' tầng tích chập, 'F' tầng kết nối đầy đủ .....	125
<b>Hình 6.5:</b> Ví dụ minh họa mạng tích chập cho nhận dạng ký tự số viết tay (nguồn [O'Neill, 2006]) .....	126
<b>Hình 6.6:</b> Kết quả bước tích chập ảnh 29x29 sử dụng mặt nạ 5x5.....	127
<b>Hình 6.7:</b> Kết quả phân lớp theo số epoch với tốc độ học 0.001 và lỗi trung bình là 0.1 ..	129
<b>Hình 6.8:</b> Kết quả phân lớp của SVM theo tham số $\gamma$ của hàm nhân RBF.....	130
<b>Hình 6.9:</b> Máy học phân lớp văn bản vào chủ đề.....	130

## DANH MỤC BẢNG

<b>Bảng 1.1: Tập dữ liệu weather</b> .....	4
<b>Bảng 2.1: Các hàm mạng thông dụng</b> .....	22
<b>Bảng 2.2: Các hàm kích hoạt dùng cho nơ-ron</b> .....	23
<b>Bảng 2.3: Giải thuật huấn luyện Perceptron</b> .....	29
<b>Bảng 2.4: Bảng chân trị của hàm AND</b> .....	30
<b>Bảng 2.5: Giải thuật huấn luyện Perceptron với luật gradient chuẩn</b> .....	36
<b>Bảng 2.6: Giải thuật huấn luyện Perceptron với luật Delta</b> .....	37
<b>Bảng 3.1: Phân lớp với máy học SVM</b> .....	62
<b>Bảng 4.1: Ma trận phân lớp C</b> .....	76
<b>Bảng 4.2: Ma trận phân lớp trình bày kết quả dự đoán nhãn của mô hình M trên tập kiểm tra TEST-DATA</b> .....	77
<b>Bảng 4.3: Ma trận phân lớp 2x2</b> .....	78
<b>Bảng 4.4: Ma trận phân lớp thu được từ mô hình M1</b> .....	78
<b>Bảng 4.5: Ma trận phân lớp thu được từ mô hình M2</b> .....	79
<b>Bảng 5.1: Giải thuật <math>k</math>-means</b> .....	85
<b>Bảng 5.2: Tập dữ liệu mô phỏng</b> .....	88
<b>Bảng 5.3: Tập dữ liệu Iris</b> .....	93
<b>Bảng 5.4: Giải thuật huấn luyện bản đồ tự tổ chức</b> .....	97
<b>Bảng 5.5: Giải thuật EM cho bài toán gom nhóm dữ liệu</b> .....	110
<b>Bảng 5.6: Giải thuật hỗn hợp các phân phối chuẩn</b> .....	116
<b>Bảng 5.7: Một số công thức liên quan đến đạo hàm trên ma trận, vec-tơ</b> .....	117
<b>Bảng 6.1: Kết quả phân lớp tập dữ liệu 20-newsgroups</b> .....	133
<b>Bảng 6.2: Kết quả phân lớp tập dữ liệu Reuters</b> .....	134

## THUẬT NGỮ

**Machine learning:** máy học

**Supervised learning:** học có giám sát

**Unsupervised learning:** học không giám sát

**Classification:** phân lớp (học có giám sát)

**Regression:** hồi quy

**Clustering:** gom cụm, gom nhóm (học không giám sát)

**Minimizing a risk functional:** cực tiểu rủi ro

**Empirical risk minimization:** cực tiểu rủi ro thực nghiệm

**Empirical inference:** suy diễn thực nghiệm

**Vapnik-Chervonenkis dimension:** chiều VC

**Structural risk minimization induction principle:** Nguyên lý quy nạp cực tiểu rủi ro cấu trúc

**Neural network:** mạng nơ-ron nhân tạo

**MLP (multilayer perceptron):** mạng perceptron đa tầng

**Convolutional neural network:** mạng nơ-ron tích chập

**Overfitting:** học vẹt

**Underfitting:** học không thuộc bài

**SVM (Support vector machines):** máy học vectơ hỗ trợ

**Hyperplane:** siêu phẳng

**Kernel:** hàm nhân

**Feature space:** không gian trung gian hay không gian đặc trưng

**Density function estimation:** ước lượng hàm mật độ xác suất

**SOM (Self-organizing map):** bản đồ tự tổ chức

**MLE (maximum likelihood estimation):** phương pháp hợp lý cực đại

**EM (Expectation maximization):** cực đại hóa kỳ vọng

**Decision tree:** cây quyết định

**Naive Bayes:** Bayes thơ ngây

**kNN (k nearest neighbors):** k láng giềng

Mẫu tin còn được gọi là điểm dữ liệu, phần tử. Thuộc tính dữ liệu cũng là trường dữ liệu, biến hay chiều dữ liệu. Một lớp của dữ liệu thường được gán một nhãn hay còn gọi là nhãn.

# **PHẦN I**

## **GIỚI THIỆU**

- *KHÁI NIỆM MÁY HỌC*
- *NGUYÊN LÝ HỌC THỐNG KÊ*

## Chương 1

# GIỚI THIỆU

Máy học, có tài liệu gọi là học máy (machine learning), là một lĩnh vực của trí tuệ nhân tạo, liên quan đến việc phát triển các kỹ thuật cho phép các máy tính có thể học từ dữ liệu. Máy học là tiếp cận để tạo ra các chương trình máy tính có thể học tự động từ dữ liệu để nhận dạng các mẫu phức tạp và đưa ra các quyết định thông minh. Vấn đề khó khăn có thể thấy ở đây là đầu vào của tập hành vi ứng xử là rất lớn, nhưng chúng ta lại có tập dữ liệu quan sát để máy học là hữu hạn, không thể phủ được toàn bộ tập hành vi ứng xử. Chương trình máy học phải tổng quát từ dữ liệu hữu hạn cho trước, có thể sinh ra hành vi ứng xử thích hợp cho những trường hợp (dữ liệu) mới đến. Máy học có tính ứng dụng rất cao trong thực tế như tìm kiếm thông tin, chẩn đoán y khoa, phát hiện thẻ tín dụng giả, phân tích thị trường chứng khoán, phân loại các chuỗi DNA, nhận dạng tiếng nói và chữ viết, dịch tự động, chơi trò chơi và điều khiển robot.

Trong chương này, chúng tôi trình bày các khái niệm cơ bản trong máy học. Trước tiên, chúng tôi giới thiệu ngắn gọn về máy học. Tiếp đến, các vấn đề chính trong máy học bao gồm học có giám sát (supervised learning, ví dụ như phân lớp, hồi quy) và học không giám sát (unsupervised learning, ví dụ như ước lượng mật độ, gom nhóm dữ liệu). Chúng tôi sẽ trình bày nguyên lý học thống kê [Vapnik, 1995], [Vapnik, 1999], là nền tảng nguyên lý học máy, nguyên lý cực tiểu rủi ro thực nghiệm, yếu tố quan trọng trong học thống kê, chiều VC (Vapnik-Chervonenkis dimension), tại sao mạng nơ-ron và máy học SVM là giải thuật học tổng quát nhất.

### 1.1 KHÁI NIỆM CƠ BẢN MÁY HỌC

Máy học là chương trình máy tính cho phép học tự động từ dữ liệu để nhận dạng các mẫu phức tạp, tạo ra hành vi ứng xử thông minh với trường hợp mới đến. [Mitchell, 1997] định nghĩa về học máy như sau: một chương trình máy tính được gọi là học từ kinh nghiệm  $E$  với một vài lớp của vấn đề  $T$  và độ đo hiệu quả  $P$ , nếu hiệu năng của vấn đề trong  $T$ , đánh giá theo tiêu chí  $P$ , được cải thiện từ kinh nghiệm  $E$ .

Kinh nghiệm  $E$  chính là dữ liệu thực nghiệm (hay còn gọi là ví dụ). Lớp vấn đề trong máy học có thể phân thành hai nhóm chính là học có giám sát (supervised learning) và học không giám sát (unsupervised learning).

Chúng ta xét một ví dụ sau đây về máy học có giám sát. Giả sử người ta cần xây dựng mô hình học tự động từ dữ liệu, để có thể dự báo chơi golf hay không. Theo đó, dữ liệu weather [Mitchell, 1997], dùng cho máy học (như bảng 1.1), có các thuộc tính outlook (quang cảnh), temperature (nhiệt độ), humidity (độ ẩm), windy (gió) và Play (lớp dự báo: chơi golf hay không chơi cho từng trường hợp cụ thể). Một giải thuật học  $A$  dùng tập dữ liệu weather để học mô

hình dự báo  $H$ , sao cho một ngày nào đó có thời tiết  $[outlook=rainy, temp=73, humidity=90, windy=FALSE]$  thì mô hình  $H$  có thể dự báo được có chơi golf hay không. Hơn nữa, giải thuật  $A$  có thể cải tiến mô hình  $H$  dựa trên tập dữ liệu học để có thể cho kết quả dự báo tốt.

**Bảng 1.1: Tập dữ liệu weather**

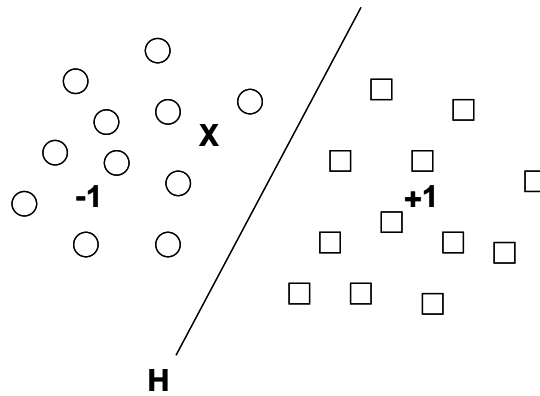
Outlook	Temp.	Humidity	Windy	Play
sunny	85	85	FALSE	no
sunny	80	90	TRUE	no
overcast	83	86	FALSE	yes
rainy	70	96	FALSE	yes
rainy	68	80	FALSE	yes
rainy	65	70	TRUE	no
overcast	64	65	TRUE	yes
sunny	72	95	FALSE	no
sunny	69	70	FALSE	yes
rainy	75	80	FALSE	Yes
sunny	75	70	TRUE	Yes
overcast	72	90	TRUE	Yes
overcast	81	75	FALSE	Yes
rainy	71	91	TRUE	no

Trong giải thích các khái niệm máy học tiếp theo sau đây, chúng ta thống nhất dùng thuật ngữ tập dữ liệu có  $m$  phần tử  $\{x_1, x_2, \dots, x_m\}$ , với  $x_i \in R^n$ , là các véc-tơ  $n$  biến (thuộc tính, chiều) ngẫu nhiên, độc lập, phân phối đồng nhất (independent and identically distributed, viết tắt là i.i.d.).

**Học có giám sát:** thuật toán học tạo ra một hàm ánh xạ dữ liệu đầu vào tới kết quả đích mong muốn (nhãn, lớp, giá trị cần dự báo). Trong học có giám sát, tập dữ liệu dùng để huấn luyện phải được gán nhãn, lớp hay giá trị cần dự báo.

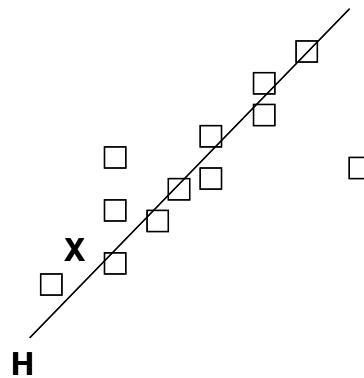
Xét bài toán phân lớp (classification) xây dựng mô hình  $H$  được huấn luyện từ tập dữ liệu  $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$  (với  $x_i \in R^n$ , có nhãn (lớp) tương ứng  $y_i$  rời rạc, ví dụ  $y_i \in \{\pm 1\}$ ), tách lớp, để khi phần tử  $x$  mới đến, mô hình  $H$  có thể gán nhãn (phân lớp)  $x$  vào một trong số lớp cho trước (ví dụ như nhãn của  $x$  được  $H$  gán là  $-1$ , xem hình 1.1).

Vấn đề học từ tập dữ liệu weather để dự báo chơi golf hay không là một ví dụ về bài toán phân lớp (trong trường hợp này,  $y_i$  là rời rạc, 2 lớp là play = yes và play = no).



**Hình 1.1:** Bài toán phân lớp

Bài toán hồi quy (regression), tương tự như bài toán phân lớp, xây dựng mô hình  $H$  được huấn luyện từ tập dữ liệu  $\{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$  (với  $x_i \in \mathbb{R}^n$ , tương ứng  $y_i$  là giá trị dự báo là liên tục,  $y_i \in \mathbb{R}^l$ ), đi qua dữ liệu sao cho sai khác là nhỏ nhất, để khi phần tử  $x$  mới đến, mô hình  $H$  có thể dự báo  $x$  là một giá trị (ví dụ như  $x$  được  $H$  giá trị dự báo là khoảng cách từ  $x$  đến  $H$  như xem hình 1.2).

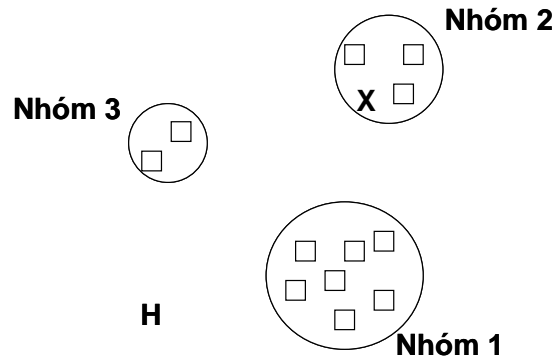


**Hình 1.2:** Bài toán hồi quy

Để minh họa cho bài toán hồi quy, chúng ta có thể xét vấn đề dự báo lượng mưa của tiểu vùng sông Mê Kông. Với tập dữ liệu thu thập bao gồm các biến như nhiệt độ, bức xạ mặt trời, độ ẩm, hướng gió, tốc độ gió, lượng mưa. Giải thuật học  $A$  cần học mô hình dự báo  $H$  để có thể dự báo lượng mưa (trong trường hợp này,  $y_i$  là số liên tục do lượng mưa cần dự báo là bao nhiêu mm) của ngày nào đó khi biết được nhiệt độ, bức xạ mặt trời, độ ẩm, hướng gió và tốc độ gió của ngày đó.



**Học không giám sát:** thuật toán học thực hiện mô hình hóa một tập dữ liệu đầu vào, không được gán nhãn (lớp, giá trị cần dự báo). Xét bài toán gom nhóm dữ liệu, xây dựng mô hình  $H$  được huấn luyện từ tập dữ liệu  $\{x_1, x_2, \dots, x_m\}$  (với  $x_i \in R^n$ ), sao cho hai phần tử dữ liệu của cùng nhóm phải có tính chất rất giống nhau hơn là hai phần tử thuộc hai nhóm khác nhau. Ví dụ như hình 1.3, xây dựng mô hình  $H$  gom nhóm dữ liệu thành ba nhóm.



**Hình 1.3:** Bài toán gom nhóm dữ liệu

Ví dụ như chúng ta cần gom nhóm sinh viên có các đặc trưng tương tự nhau như sở thích, kết quả học tập, trình độ ngoại ngữ, tính năng động.

Ngoài ra, còn có lớp bài toán như:

**Học bán giám sát** (semi-supervised learning): kết hợp dữ liệu có gán nhãn và không gán nhãn để tạo một bộ phân lớp.

**Học tăng cường** (reinforcement learning): thuật toán học một ứng xử dựa vào quan sát về thế giới. Mỗi hành động đều có tác động tới môi trường và môi trường cung cấp thông tin phản hồi để hướng dẫn cho thuật toán trong quá trình học.

## 1.2 NGUYÊN LÝ HỌC THỐNG KÊ

Để nghiên cứu nền tảng lý thuyết của nguyên lý máy học, chúng tôi trình bày nguyên lý học thống kê [Vapnik, 1995], [Vapnik, 1999]. Trước tiên, chúng ta sẽ làm quen với mô hình ước lượng hàm, cực tiểu hóa rủi ro, đây là khung tổng quát của các bài toán học máy như phân lớp, hồi quy, ước lượng mật độ dữ liệu.

### 1.2.1 Mô hình ước lượng hàm

Mô hình máy học từ dữ liệu có thể được mô tả bởi ba bộ phận như sau.

1. Bộ sinh các véc-tơ ngẫu nhiên  $x$ , độc lập, cùng phân phối xác suất cố định nhưng chưa được biết trước  $P(x)$ .
2. Một giám sát (supervisor) trả về ở đầu ra một véc-tơ  $y$  khi nhận ở đầu vào véc-tơ  $x$ , theo hàm phân phối xác suất có điều kiện  $P(y|x)$  cố định nhưng chưa được biết trước.
3. Một máy học có thể cài đặt một tập hàm  $f(x, \alpha)$ , với  $\alpha \in A$   
*( $A$  là tập hợp các tham số,  $\alpha$  không nhất thiết là một véc-tơ, chính vì vậy  $f(x, \alpha)$  được dùng ở đây là tập hàm bất kỳ).*

Vấn đề học máy chính là tìm kiếm (chọn lựa) từ tập hàm cho trước  $f(x, \alpha)$ ,  $\alpha \in A$ , một hàm mà sao cho kết quả dự báo thu được gần nhất với kết quả của giám sát. Sự lựa chọn hàm dựa vào tập dữ liệu huấn luyện có  $m$  phần tử:

$$(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m) \tag{1.1}$$

với  $x_i \in R^n$ , là các véc-tơ  $n$  biến (thuộc tính, chiều), có nhãn (lớp) tương ứng  $y_i$ , và  $P(x, y) = P(x)P(y|x)$ .

### 1.2.2 Cực tiểu hóa rủi ro (minimizing a risk functional)

Xem xét vấn đề lựa chọn hàm xấp xỉ tốt nhất cho dự báo của giám sát, cũng chính là độ đo sự sai khác (loss)  $L(y, f(x, \alpha))$  giữa giá trị thực, cần được dự báo  $y$  của giám sát (kết quả mong đợi, đích) và giá trị dự báo  $f(x, \alpha)$  thu được từ mô hình máy học khi nhận đầu vào  $x$ . Sự sai khác được tính thông qua rủi ro  $R(\alpha)$  như sau:

$$R[f] = \int L(y, f(x, \alpha)) dP(x, y) \tag{1.2}$$

Mục tiêu của máy học là tìm hàm  $f(x, \alpha^*)$  từ lớp hàm  $f(x, \alpha)$ ,  $\alpha \in A$ , để cực tiểu rủi ro  $R[f]$  như công thức 1.2 trong khi chưa biết về  $P(x, y)$ , chỉ có được tập dữ liệu trong (1.1). Nguyên lý cực tiểu hóa rủi ro dựa vào hàm sai khác (loss) là khung tổng quát cho cả ba bài toán phân lớp, hồi quy và ước lượng mật độ dữ liệu.

#### Bài toán phân lớp:

Tập dữ liệu  $\{x, y\} = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$  (với  $x_i \in R^n$ , có nhãn (lớp) tương ứng  $y_i \in \{\pm 1\}$ ), lớp hàm  $f(x, \alpha)$  với  $\alpha \in A$ , hàm sai khác (loss):

$$L(y, f(x, \alpha)) = \begin{cases} 0 & y = f(x, \alpha) \\ 1 & y \neq f(x, \alpha) \end{cases} \tag{1.3}$$