

ĐẠI HỌC THÁI NGUYÊN
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG



TRƯỜNG ĐỨC CƯỜNG

PHÂN CỤM DỮ LIỆU SỬ DỤNG
GIẢI THUẬT DI TRUYỀN VÀ MẠNG NƠ RON

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

Thái Nguyên - 2012

ĐẠI HỌC THÁI NGUYÊN
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG



TRƯƠNG ĐỨC CƯỜNG

PHÂN CỤM DỮ LIỆU SỬ DỤNG
GIẢI THUẬT DI TRUYỀN VÀ MẠNG NƠ RON

Chuyên ngành : *Khoa học máy tính*

Mã số : *60.48.01*

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

NGƯỜI HƯỚNG DẪN KHOA HỌC

TS. Vũ Mạnh Xuân

Thái Nguyên - 2012

LỜI CẢM ƠN

Em xin bày tỏ lòng biết ơn sâu sắc tới TS. Vũ Mạnh Xuân, thầy đã hướng dẫn, chỉ dạy tận tình để em hoàn thành luận văn này. Em xin chân thành cảm ơn các thầy, cô giáo Trường Đại học Công nghệ Thông tin & Truyền thông - Đại học Thái Nguyên, cùng các thầy, cô giáo Viện Công nghệ Thông tin - Viện Khoa học và Công nghệ Việt Nam đã truyền thụ kiến thức cho em trong suốt quá trình học tập vừa qua.

Tôi cũng xin cảm ơn cơ quan, bạn bè đồng nghiệp, gia đình và những người thân đã cùng chia sẻ, giúp đỡ, động viên, tạo mọi điều kiện thuận lợi để tôi có thể học tập và hoàn thành bản luận văn này.

Tuy đã có những cố gắng nhất định nhưng do thời gian và trình độ có hạn nên chắc chắn luận văn còn nhiều thiếu sót và hạn chế nhất định. Rất mong nhận được sự góp ý của thầy cô và các bạn.

Thái Nguyên, ngày 27 tháng 06 năm 2012

Học viên

Trương Đức Cường

LỜI CAM ĐOAN

Tôi xin cam đoan kết quả đạt được trong luận văn là sản phẩm nghiên cứu, tìm hiểu của riêng cá nhân tôi. Trong toàn bộ nội dung của luận văn, những điều được trình bày hoặc là của cá nhân tôi hoặc là được tổng hợp từ nhiều nguồn tài liệu. Tất cả các tài liệu tham khảo đều có xuất xứ rõ ràng và được trích dẫn hợp pháp.

Tôi xin hoàn toàn chịu trách nhiệm và chịu mọi hình thức kỷ luật theo quy định cho lời cam đoan của mình.

Thái Nguyên, ngày 27 tháng 06 năm 2012

Học viên

Trương Đức Cường

MỤC LỤC

MỤC LỤC	iii
DANH SÁCH HÌNH VẼ	v
DANH SÁCH BẢNG BIỂU	vi
DANH SÁCH TỪ VIẾT TẮT	vii
MỞ ĐẦU	1
CHƯƠNG I: TỔNG QUAN VỀ PHÂN CỤM DỮ LIỆU	3
<i>1.1. Khái niệm và mục đích của phân cụm dữ liệu</i>	<i>3</i>
<i>1.2. Ứng dụng của phân cụm dữ liệu</i>	<i>4</i>
<i>1.3. Một số phương pháp phân cụm dữ liệu</i>	<i>5</i>
1.3.1. Phân cụm phân hoạch	<i>5</i>
1.3.2. Phân cụm phân cấp	<i>7</i>
1.3.3. Phân cụm dựa trên mật độ.....	<i>9</i>
1.3.4. Phân cụm dựa trên lưới	<i>11</i>
1.3.5. Phân cụm dữ liệu dựa trên mô hình	<i>13</i>
1.3.6. Phân cụm dữ liệu mờ	<i>14</i>
CHƯƠNG II: PHÂN CỤM DỮ LIỆU SỬ DỤNG GIẢI THUẬT DI TRUYỀN VÀ MẠNG NƠI RON	16
<i>2.1. Giải thuật di truyền</i>	<i>16</i>
2.1.1. Sơ đồ thực hiện giải thuật di truyền.....	<i>17</i>
2.1.2. Các quá trình chính trong giải thuật di truyền.....	<i>19</i>

2.1.2.1. Biểu diễn các cá thể	19
2.1.2.2. Hàm mục tiêu (Fitness).....	21
2.1.2.3. Toán tử tái tạo (Reproduction).....	21
2.1.2.4. Toán tử lai ghép (Crossover)	24
2.1.2.5. Toán tử đột biến (Mutation).....	26
2.1.2.6. Các thông số cơ bản của giải thuật di truyền	27
2.1.3. Ưu và nhược điểm của giải thuật di truyền	28
2.2. <i>Mạng nơ ron</i>	30
2.2.1. Định nghĩa	30
2.2.2. Nơ ron sinh học và mạng nơ ron sinh học.....	31
2.2.3. Cấu trúc mạng nơ ron	32
2.2.4. Phân loại mạng nơ ron.....	33
2.3. <i>Mối quan hệ giữa giải thuật di truyền và mạng nơ ron trong phân cụm dữ liệu</i>	35
2.3.1. Một số phương thức kết hợp giữa GA và mạng nơ ron	36
2.3.2. Một số ví dụ về việc kết hợp giữa GA và mạng nơ ron	38
CHƯƠNG III: BÀI TOÁN ỨNG DỤNG.....	42
3.1. <i>Phát biểu bài toán</i>	42
3.2. <i>Thuật toán</i>	42
3.3. <i>Kết quả thử nghiệm</i>	48
3.4. <i>Nhận xét</i>	49
KẾT LUẬN	54
TÀI LIỆU THAM KHẢO	55

DANH SÁCH HÌNH VẼ

Hình 1.1. Quy trình phân cụm.....	3
Hình 1.2. Mô phỏng sự phân cụm dữ liệu	4
Hình 1.3. Các chiến lược phân cụm phân cấp.....	8
Hình 1.4. Một số hình dạng khám phá bởi phân cụm dựa trên mật độ	10
Hình 1.5. Mô hình cấu trúc dữ liệu lưới.....	12
Hình 2.1. Lưu đồ giải thuật di truyền.....	18
Hình 2.2. Bánh xe trọng số.....	23
Hình 2.3. Lai ghép một điểm	25
Hình 2.4. Lai ghép trong biểu diễn bằng giá trị	26
Hình 2.5. Cấu tạo của nơ ron	31
Hình 2.6. Thu nhận tín hiệu trong nơ ron	31
Hình 2.7. Mạng nơ ron truyền thẳng nhiều lớp.....	34
Hình 2.8. Mạng hồi quy (Recurrent Neural Network).....	34
Hình 2.9. Mô đun ghép cặp Di truyền – Nơ ron trong một hệ thống ứng dụng	38
Hình 2.10. Sơ đồ của hệ thống <i>XROUTE</i> (Kadaba, Nygard và Juell 1991) ...	38
Hình 3.1. Dữ liệu đầu ra	42
Hình 3.2. Dữ liệu đầu vào sau khi mã hóa.....	43
Hình 3.3. Quá trình lai ghép.....	43
Hình 3.4. Tập điểm dữ liệu vào.....	48
Hình 3.5. Giao diện chương trình	49
Hình 3.6. Kết quả phân cụm với string count = 100.....	50
Hình 3.7. Kết quả phân cụm với string count = 1	50
Hình 3.8. Kết quả phân cụm bộ dữ liệu giao nhau với stringcount = 1	51

DANH SÁCH BẢNG BIỂU

Bảng 2.1. Bảng thể hiện tổng giá trị hàm mục tiêu.....	22
Bảng 2.2. Chuỗi nhiễm sắc thể	23
Bảng 2.3. Lai ghép mặt nạ	25
Bảng 2.4. So sánh <i>K-mean</i> và Giải thuật di truyền.....	29
Bảng 2.5. Một số phương thức liên GA với mạng nơ ron	37

DANH SÁCH TỪ VIẾT TẮT

Từ viết tắt	Ý nghĩa
KPDL	Khai phá dữ liệu
PCDL	Phân cụm dữ liệu
CSDL	Cơ sở dữ liệu
GA	Giải thuật di truyền
NST	Nhiễm sắc thể

MỞ ĐẦU

Trong những năm gần đây, sự phát triển mạnh mẽ của công nghệ thông tin và ngành công nghiệp phần cứng đã làm cho khả năng thu thập và lưu trữ thông tin của các hệ thống thông tin tăng nhanh một cách chóng mặt. Bên cạnh đó việc tin học hóa một cách ồ ạt và nhanh chóng các hoạt động sản xuất, kinh doanh cũng như nhiều lĩnh vực hoạt động khác đã tạo cho chúng ta một hệ thống cơ sở dữ liệu khổng lồ. Hệ thống này đã đem lại những lợi ích vô cùng to lớn cho con người trong việc lưu trữ, tìm kiếm và thống kê. Tuy vậy, sự bùng nổ này đã dẫn tới một nhu cầu mới là phát hiện tri thức từ kho dữ liệu khổng lồ đó. Đây là một vấn đề rất phức tạp, cần phải có những công cụ và kỹ thuật xử lý linh hoạt như suy nghĩ của con người.

Trong ngành khoa học máy tính, tìm kiếm lời giải tối ưu cho các bài toán là vấn đề được các nhà khoa học máy tính đặc biệt rất quan tâm. Mục đích chính của các thuật toán là tìm kiếm thuật giải chất lượng cao và sử dụng kỹ thuật trí tuệ nhân tạo đặc biệt rất cần thiết khi giải quyết các bài toán có không gian tìm kiếm lớn.

Giải thuật di truyền (Genetic Algorithm - GA) là một trong những kỹ thuật tìm kiếm lời giải tối ưu đã đáp ứng được yêu cầu của nhiều bài toán và ứng dụng. Hiện nay, thuật toán di truyền cùng với mạng nơ ron được ứng dụng rất rộng rãi trong các lĩnh vực phức tạp. Thuật toán di truyền kết hợp với mạng nơ ron chứng tỏ được hiệu quả của nó trong các vấn đề khó có thể giải quyết bằng các phương pháp thông thường hay các phương pháp cổ điển, nhất là trong các bài toán cần có sự lượng giá, đánh giá sự tối ưu của kết quả thu được.

Chính vì vậy, trong phạm vi đề tài này, tôi chọn hướng kết hợp giữa mạng nơ ron và giải thuật di truyền áp dụng vào bài toán phân cụm dữ liệu,