

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



Nguyễn Như Thế

**NGHIÊN CỨU CÁC PHƯƠNG PHÁP PHÂN LỚP DỮ LIỆU
VÀ ỨNG DỤNG TRONG BÀI TOÁN DỰ BÁO THUÊ BAO
RỜI MẠNG VIỄN THÔNG**

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

Thái Nguyên - 2016

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



Nguyễn Như Thế

**NGHIÊN CỨU CÁC PHƯƠNG PHÁP PHÂN LỚP DỮ LIỆU
VÀ ỨNG DỤNG TRONG BÀI TOÁN DỰ BÁO THUÊ BAO
RỜI MẠNG VIỄN THÔNG**

Chuyên ngành: Khoa học máy tính

Mã số: 60 48 0101

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

NGƯỜI HƯỚNG DẪN KHOA HỌC: TS NGUYỄN LONG GIANG

Thái Nguyên - 2016

LỜI CAM ĐOAN

Tên tôi là: **Nguyễn Như Thế**

Sinh ngày: 12/12/1989

Học viên lớp cao học: CHK13E - Trường Đại học Công nghệ thông tin và Truyền thông – Đại học Thái Nguyên.

Hiện đang công tác tại: Sở Thông tin và Truyền thông tỉnh Phú Thọ

Xin cam đoan: Đề tài *“Nghiên cứu các phương pháp phân lớp dữ liệu và ứng dụng trong bài toán dự báo thuê bao rờì mạng viễn thông”* do Thầy giáo **TS. Nguyễn Long Giang** hướng dẫn là công trình nghiên cứu của riêng tôi. Tất cả tài liệu tham khảo đều có nguồn gốc, xuất xứ rõ ràng.

Tác giả xin cam đoan tất cả những nội dung trong luận văn đúng như nội dung trong đề cương và yêu cầu của thầy giáo hướng dẫn. Nếu sai tôi hoàn toàn chịu trách nhiệm trước hội đồng khoa học và trước pháp luật.

Thái Nguyên, ngày 28 tháng 6 năm 2016

HỌC VIÊN

Nguyễn Như Thế

LỜI CẢM ƠN

Sau một thời gian nghiên cứu và làm việc nghiêm túc, được sự động viên, giúp đỡ và hướng dẫn tận tình của Thầy giáo hướng dẫn **TS. Nguyễn Long Giang**, luận văn với đề tài “*Nghiên cứu các phương pháp phân lớp dữ liệu và ứng dụng trong bài toán dự báo thuê bao rời mạng viễn thông*” đã hoàn thành.

Tôi xin bày tỏ lòng biết ơn sâu sắc đến:

Thầy giáo hướng dẫn **TS. Nguyễn Long Giang** đã tận tình chỉ dẫn, giúp đỡ tôi hoàn thành luận văn này.

Tôi xin bày tỏ lòng biết ơn đến các thầy cô trong Trường Đại học Công nghệ thông tin và Truyền thông – Đại học Thái Nguyên đã giúp đỡ tôi trong quá trình học tập cũng như thực hiện luận văn.

Tôi xin cảm ơn Chi nhánh Mobifone Phú Thọ đã nhiệt tình giúp đỡ, cung cấp thông tin trong quá trình nghiên cứu, thực nghiệm chương trình luận văn.

Tôi xin chân thành cảm ơn bạn bè, đồng nghiệp và gia đình đã động viên, khích lệ, tạo điều kiện giúp đỡ tôi trong suốt quá trình học tập, thực hiện và hoàn thành luận văn này.

Thái Nguyên, ngày 28 tháng 6 năm 2016

HỌC VIÊN

Nguyễn Như Thế

MỤC LỤC

| | |
|--|-----|
| LỜI CAM ĐOAN | i |
| LỜI CẢM ƠN | ii |
| DANH MỤC CÁC KÝ HIỆU VÀ CHỮ VIẾT TẮT | v |
| DANH MỤC HÌNH ẢNH..... | vi |
| DANH MỤC BẢNG BIỂU..... | vii |
| MỞ ĐẦU | 1 |
| <i>Chương 1: TỔNG QUAN VỀ KHAI PHÁ DỮ LIỆU</i> | 3 |
| 1.1. Tổng quan về khai phá dữ liệu..... | 3 |
| 1.1.1. Tại sao cần khai phá dữ liệu | 3 |
| 1.1.2. Các khái niệm cơ bản | 3 |
| 1.1.3. Quy trình khai phá dữ liệu..... | 5 |
| 1.1.4. Các bài toán cơ bản trong khai phá dữ liệu | 6 |
| 1.1.5. Các ứng dụng của khai phá dữ liệu..... | 7 |
| 1.1.6. Quy trình xây dựng mô hình khai phá dữ liệu..... | 8 |
| 1.2. Bài toán phân lớp và dự báo | 10 |
| 1.2.1. Giới thiệu bài toán..... | 10 |
| 1.2.2. Các bước giải quyết bài toán | 11 |
| <i>Chương 2: CÁC PHƯƠNG PHÁP PHÂN LỚP TRONG KHAI PHÁ DỮ LIỆU</i> | 12 |
| 2.1. Phân lớp bằng phương pháp quy nạp cây quyết định | 12 |
| 2.2. Phân lớp bằng phương pháp Bayesian | 15 |
| 2.3. Support Vector Machine (SVM)..... | 16 |
| 2.3.1 Phân tách tuyến tính với lề cực đại | 16 |
| 2.3.1.1. Tìm kiếm siêu phẳng với lề cực đại | 21 |
| 2.3.1.2. Hàm phân loại tuyến tính với lề mềm cực đại..... | 22 |
| 2.3.1.3. Lý thuyết tối ưu Lagrangian | 23 |
| 2.3.1.4. Tìm kiếm siêu phẳng với lề cực đại | 25 |
| 2.3.2. Phương pháp hàm nhân (kernel methods)..... | 28 |
| 2.3.2.1 Chiều VC về khả năng phân tách của hàm tuyến tính | 29 |
| 2.3.2.2 Hàm nhân và SVM phi tuyến (Kernel function and nonlinear SVMs) . | 30 |

| | |
|--|----|
| 2.4. Phân lớp bằng mạng lan truyền ngược (mạng Nơron)..... | 33 |
| Chương 3: ỨNG DỤNG BÀI TOÁN PHÂN LỚP DỮ LIỆU THUÊ BAO RỜI MẠNG VIỄN THÔNG | 37 |
| 3.1. Bài toán phân lớp dữ liệu thuê bao rời mạng..... | 37 |
| 3.1.1. Phát biểu bài toán..... | 37 |
| 3.1.2. Khái niệm thuê bao rời mạng “churn” | 38 |
| 3.1.3. Thu thập, chuẩn hóa dữ liệu | 39 |
| 3.1.4. Lựa chọn thuộc tính..... | 42 |
| 3.2. Lựa chọn phương pháp, công cụ..... | 45 |
| 3.2.1. Ngôn ngữ R..... | 45 |
| 3.2.2. Phương pháp phân lớp..... | 47 |
| 3.2.3. Đánh giá hiệu năng..... | 48 |
| 3.3. Thực nghiệm phân lớp trên ngôn ngữ R..... | 50 |
| 3.3.1. Phân lớp dữ liệu sử dụng cây quyết định C4.5..... | 51 |
| 3.3.2. Phân lớp dữ liệu sử dụng phương pháp Naive Bayes | 53 |
| 3.3.3. Phân lớp dữ liệu bằng Support Vector Machines | 55 |
| 3.3. Đánh giá kết quả..... | 56 |
| KẾT LUẬN | 58 |
| TÀI LIỆU THAM KHẢO..... | 60 |

DANH MỤC CÁC KÝ HIỆU VÀ CHỮ VIẾT TẮT

| TT | Thuật ngữ | Định nghĩa |
|-----------|------------------|-------------------------------------|
| 1. | KPDL | Khai phá dữ liệu |
| 2. | KDD | Knowledge Discovery and Data Mining |
| 3. | NB | Naïve Bayes |
| 4. | SVM | Support vector machine |
| 5. | NN | Neural Networks |

DANH MỤC HÌNH ẢNH

| | |
|--|----|
| Hình 1.1- Các bước trong khai phá dữ liệu | 6 |
| Hình 1.2 - Quy trình xây dựng mô hình khai phá dữ liệu | 9 |
| Hình 2.1 - Ví dụ về cây quyết định | 12 |
| Hình 2.2 - Về mặt trực quan thì hàm tuyến tính siêu phẳng với lẽ lớn nhất trông có vẻ hợp lý | 19 |
| Hình 2.3 - Ví dụ về bài toán phân loại trong không gian hai chiều | 19 |
| Hình 2.4 - Ba điểm trong mặt phẳng bị chia tách bởi một đường thẳng có hướng. .. | 28 |
| Hình 2.5 - Mạng nơ-ron truyền thẳng nhiều lớp | 34 |
| Hình 3.1 - Mô hình quan hệ các bảng dữ liệu | 40 |
| Hình 3.2 - Các giai đoạn của mô hình dự đoán thuê bao rời mạng | 42 |
| Hình 3.3 - Lựa chọn thuộc tính trong phân lớp dữ liệu | 44 |
| Hình 3.4 - Số lượng thuộc tính được thu thập | 44 |
| Hình 3.5 – Giao diện làm việc trên ngôn ngữ R | 47 |
| Hình 3.6 – Quy trình thực nghiệm bài toán phân lớp dữ liệu thuê bao rời mạng | 50 |
| Hình 3.7- mô hình phân lớp cây quyết định | 52 |
| Hình 3.8 - Chi tiết nút nhánh thứ 15 trong phân lớp cây quyết định | 52 |
| Hình 3.9 - Kết quả phân lớp dữ liệu bằng SVM | 55 |
| Hình 3.10 – Hiệu năng các thuật toán với lớp thuê bao rời mạng | 57 |

DANH MỤC BẢNG BIỂU

| | |
|---|----|
| Bảng 1 - Ma trận nhầm lẫn | 49 |
| Bảng 2 – Kết quả mô hình phân lớp sử dụng C 4.5 | 53 |
| Bảng 3 – Độ đo hiệu năng thuật toán Cây quyết định | 53 |
| Bảng 4 – Kết quả mô hình phân lớp sử dụng NB | 54 |
| Bảng 5. – Độ đo hiệu năng thuật toán NB..... | 54 |
| Bảng 6 – Kết quả mô hình phân lớp sử dụng SVM | 55 |
| Bảng 7. – Độ đo hiệu năng thuật toán SVM..... | 56 |
| Bảng 8. – Tổng hợp đánh giá hiệu năng các phương pháp phân lớp..... | 56 |

MỞ ĐẦU

Sự bùng nổ và phát triển của ngành công nghệ thông tin đã làm lượng dữ liệu được thu thập và lưu trữ ở các hệ thống thông tin tăng lên một cách nhanh chóng. Trước tình hình đó, việc khai thác và chọn lọc những dữ liệu có ích, tiềm ẩn từ lượng dữ liệu khổng lồ này là rất cần thiết. Các tri thức trích lọc từ dữ liệu sẽ giúp các cơ quan, tổ chức đưa ra những dự báo và điều hành hiệu quả.

Khai phá dữ liệu và khám phá tri thức (Data mining and Knowledge discovery) là một lĩnh vực quan trọng của ngành Công nghệ thông tin với mục tiêu là tìm kiếm các tri thức có ích, cần thiết, tiềm ẩn và chưa được biết trước trong cơ sở dữ liệu lớn. Đây là lĩnh vực đã và đang thu hút đông đảo các nhà khoa học trên thế giới và trong nước tham gia nghiên cứu. Phân lớp (classification) là một trong những bài toán cơ bản trong khai phá dữ liệu với mục tiêu là phân loại các đối tượng vào các lớp cho trước. Theo tiếp cận học máy, phân lớp là phương pháp học có giám sát với hai giai đoạn: Giai đoạn 1 là xây dựng mô hình phân lớp dựa trên tập dữ liệu huấn luyện có đầu vào và đầu ra mong muốn (gọi là nhãn lớp); Giai đoạn 2 là sử dụng mô hình phân lớp để phân loại các tập dữ liệu chưa có nhãn lớp vào các lớp đã cho và có ứng dụng trong nhiều bài toán dự báo trong thực tế. Phân lớp được sử dụng rộng rãi trong các bài toán thực tiễn như trong y tế, ngân hàng, viễn thông, kinh tế, tài chính...

Ngày nay, cùng với sự phát triển mạnh mẽ của thị trường viễn thông là sự ra đời của nhiều nhà cung cấp và kinh doanh dịch vụ mạng viễn thông. Thị trường viễn thông đang đi vào giai đoạn bão hòa, khách hàng có nhiều sự lựa chọn, dẫn đến họ có thể thay đổi sử dụng dịch vụ bất cứ khi nào, kết quả là số