

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

ĐÀO MỸ HẠNH

CỤM DỮ LIỆU VÀ ỨNG DỤNG TRONG PHÂN TÍCH LƯƠNG CỦA CÁN BỘ

TRƯỜNG CAO ĐẲNG NGHỀ HÀ NAM

Chuyên ngành: KHOA HỌC MÁY TÍNH

Mã số chuyên ngành: 60 48 0101

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

Thái Nguyên - 2015

LỜI CẢM ƠN

Tôi xin chân thành cảm ơn tập thể các thầy cô trong khoa đào tạo sau đại học trường Đại học Công nghệ Thông tin và Truyền thông Thái Nguyên đã trang bị cho tôi những kiến thức cơ bản trong những năm học tập tại trường để tôi có thể hoàn thành tốt bản luận văn tốt nghiệp này.

Tôi xin cảm ơn các đồng nghiệp và người thân đã động viên, giúp đỡ tôi trong quá trình nghiên cứu và thực hiện luận văn.

Đặc biệt, tôi xin cảm ơn **GS.TS Vũ Đức Thi**, người đã trực tiếp, tận tâm hướng dẫn, giúp đỡ, cung cấp tài liệu và tạo mọi điều kiện thuận lợi cho tôi nghiên cứu thành công luận văn tốt nghiệp của mình.

Thái Nguyên, ngày ... tháng ... năm 2015

Tác giả luận văn

Đào Mỹ Hạnh

LỜI CAM ĐOAN

Tôi xin cam đoan toàn bộ nội dung bản luận văn này là do tôi tự sưu tầm, tra cứu và sắp xếp cho phù hợp với nội dung yêu cầu của đề tài.

Nội dung luận văn này chưa từng được công bố hay xuất bản dưới bất kỳ hình thức nào và cũng không được sao chép từ bất kỳ một công trình nghiên cứu nào.

Các số liệu, kết quả nêu trong luận văn là trung thực và chưa từng được ai công bố trong bất kỳ công trình nào khác. Tôi cũng xin cam đoan rằng mọi sự giúp đỡ cho việc thực hiện luận văn này đã được cảm ơn và các thông tin trích dẫn trong luận văn đã được chỉ rõ nguồn gốc.

Nếu sai tôi xin hoàn toàn chịu trách nhiệm.

Thái Nguyên, ngày ... tháng ... năm 2015

Người cam đoan

Đào Mỹ Hạnh

DANH MỤC TỪ VIẾT TẮT

CSDL: Cơ sở dữ liệu

KPDL: Khai phá dữ liệu

PCDL: Phân cụm dữ liệu

DANH MỤC CÁC BẢNG

| | |
|--|----|
| Bảng 1.1: Thuộc tính dữ liệu nhị phân..... | 8 |
| Bảng 2. 1: Các nhóm cơ sở tương ứng..... | 43 |

DANH MỤC HÌNH VẼ

| | |
|--|----|
| Hình 1.1: Phân cụm dữ liệu..... | 5 |
| Hình 1.2: Ví dụ minh họa phân cụm phân hoạch..... | 11 |
| Hình 2.1: Kết quả phân nhóm thuật toán K–Means (a), Seed–Kmeans (b)..... | 18 |
| Hình 2.2: Lân cận của p với ngưỡng Eps..... | 18 |
| Hình 2.3: Mật độ đến được trực tiếp..... | 19 |
| Hình 2.4: Mật độ đến được..... | 19 |
| Hình 2.5: Mật độ liên thông..... | 20 |
| Hình 2.6: Đồ thị đã sắp xếp 4-dist đối với CSDL mẫu 3..... | 23 |
| Hình 2.7: Các nhóm phát hiện được bởi và DBSCAN..... | 23 |
| Hình 2.8: Các đối tượng bị ảnh hưởng trong một CSDL mẫu..... | 27 |
| Hình 2.9: Các trường hợp khác nhau của thuật toán..... | 30 |
| Hình 2.10: Thể hiện trộn các nhóm A, B, C bằng thuật toán thêm..... | 31 |

| | |
|---|----|
| Hình 2.11: Các trường hợp khác nhau của thuật toán xóa | 32 |
| Hình 2.12: Suffix trie và cây hậu tố của xâu $S = \text{abaab}$ | 35 |
| Hình 2.13: Cây hậu tố cho chuỗi $S = \text{xabxac}$ | 36 |
| Hình 2.14: Các bước tạo cây hậu tố của xâu $S = \text{abaab}$ | 37 |
| Hình 2.15: Quy tắc thêm kí tự ai vào cây đã chứa ai | 37 |
| Hình 2.16: Cây hậu tố T của xâu $S = \text{axabx}$ | 38 |
| Hình 2.17: Cây hậu tố T của xâu $S = \text{axabxb}$ theo quy tắc 1 | 38 |
| Hình 2.18: Cây hậu tố T của xâu $S = \text{axabxb}$ theo quy tắc 2 | 39 |
| Hình 2.19: Cây hậu tố với các liên kết hậu tố cho 2 chuỗi xabxa và abxbx | 40 |
| Hình 2.20: Cây hậu tố của các chuỗi "cat ate cheese", "mouse ate cheese too" and "cat ate mouse too" | 43 |
| Hình 2.21: Đồ thị các nhóm cơ sở | 44 |
| Hình 3.1: Mô hình 3-Tier. | 54 |
| Hình 3.2: Mô hình use case tổng quan hệ thống. | 55 |
| Hình 3.3: Giao diện form đăng nhập | 56 |
| Hình 3.4: Giao diện form quản lý danh mục | 57 |
| Hình 3.5: Màn hình chính | 58 |
| Hình 3.6: Dữ liệu đầu vào | 59 |
| Hình 3.7: Kết quả phân cụm dữ liệu bởi Incremental DBSCAN | 60 |
| Hình 3.8: Dữ liệu được thêm mới | 61 |
| Hình 3.9: Kết quả phân cụm sau khi thêm dữ liệu mới | 61 |
| Hình 3.10: Màn hình quản lý người dùng | 62 |
| Hình 3.11: Màn hình thêm mới người dùng | 62 |
| Hình 3.12: Màn hình sửa thông tin người dùng | 63 |
| Hình 3.13: Cửa sổ xác thực xóa thông tin người dùng | 63 |
| Hình 3.14: Màn hình quản lý thông tin khoa/viện | 64 |
| Hình 3.15: Màn hình quản lý thông tin giảng viên | 64 |
| Hình 3.16: Màn hình quản lý thông tin giảng viên | 65 |

MỤC LỤC

| | |
|---|-----|
| LỜI CẢM ƠN..... | i |
| LỜI CAM ĐOAN | iii |
| DANH MỤC TỪ VIẾT TẮT | iv |
| DANH MỤC CÁC BẢNG | iv |
| DANH MỤC HÌNH VẼ | iv |
| MỤC LỤC | vi |
| MỞ ĐẦU | ix |
| CHƯƠNG I: TỔNG QUAN VỀ KHAI PHÁ DỮ LIỆU | 1 |
| VÀ PHÂN CỤM DỮ LIỆU | 1 |
| 1.1 Khai phá dữ liệu | 1 |
| 1.1.1 Giới thiệu về khai phá dữ liệu | 1 |
| 1.1.2 Quá trình khai phá dữ liệu..... | 1 |
| 1.1.3 Các kỹ thuật khai phá dữ liệu..... | 2 |
| 1.1.4 Ứng dụng của Khai phá dữ liệu..... | 3 |
| 1.1.5 Các xu thế và vấn đề cần giải quyết trong khai phá dữ liệu..... | 3 |
| 1.2 Kỹ thuật phân cụm trong Khai phá dữ liệu | 4 |

| | |
|--|----|
| 1.2.1 Tổng quan về kỹ thuật phân cụm | 4 |
| 1.2.2 Một số khái niệm cần thiết khi tiếp cận phân cụm dữ liệu | 6 |
| 1.2.2.1 Các kiểu dữ liệu và thuộc tính trong phép phân cụm..... | 6 |
| 1.2.2.2 Đo độ tương đồng..... | 7 |
| 1.2.3 Các yêu cầu đối với kỹ thuật phân cụm dữ liệu | 9 |
| 1.2.4 Các hướng tiếp cận trong phân cụm dữ liệu | 11 |
| 1.2.4.1 Phương pháp phân hoạch: | 11 |
| 1.2.4.2 Phương pháp phân cụm phân cấp..... | 12 |
| 1.2.4.3 Phương pháp phân cụm dựa trên mật độ..... | 13 |
| 1.2.4.4 Phương pháp phân cụm dựa trên lưới | 13 |
| CHƯƠNG II: | 15 |
| MỘT SỐ THUẬT TOÁN PHÂN CỤM DỮ LIỆU ĐIỂN HÌNH | 15 |
| 2.1 Thuật toán K-Means | 15 |
| 2.2 Thuật toán DBSCAN..... | 18 |
| 2.3 Thuật toán BIRCH..... | 24 |
| 2.4 Thuật toán INCREMENTAL DBSCAN..... | 25 |
| 2.4.1 Các đối tượng bị ảnh hưởng..... | 26 |
| 2.4.2 Trường hợp thêm..... | 29 |
| 2.4.3 Trường hợp xóa | 31 |
| 2.5 Thuật toán phân nhóm cây hậu tố | 34 |
| 2.5.1 Cây hậu tố..... | 34 |
| 2.5.2 Cây hậu tố - Cây hậu tố tổng quát..... | 39 |
| 2.5.3 Thuật toán STC | 41 |
| 2.6 Thuật toán dựa vào phân loại véc-tơ hỗ trợ | 46 |
| 2.6.1 Phương pháp SVM..... | 46 |
| 2.6.2 Phương pháp FSVM..... | 48 |
| CHƯƠNG III:..... | 52 |
| ỨNG DỤNG PHƯƠNG PHÁP PHÂN NHÓM DỮ LIỆU | 52 |

| | |
|---|----|
| VÀO PHÂN TÍCH LƯƠNG CỦA CÁN BỘ | 52 |
| TRƯỜNG CAO ĐẲNG NGHỀ HÀ NAM..... | 52 |
| 3.1 Đặt vấn đề..... | 52 |
| 3.2 Giải quyết vấn đề:..... | 53 |
| 3.2.1 Công cụ lựa chọn xây dựng chương trình phần mềm : | 53 |
| 3.2.2. Biểu đồ phân cấp chức năng..... | 54 |
| 3.2.3 Mô hình tổng quan hệ thống | 55 |
| 3.2.4 Thiết kế giao diện chương trình: | 56 |
| 3.2.4.1. Giao diện form đăng nhập:..... | 56 |
| 3.2.4.2. Giao diện form quản lý danh mục:..... | 56 |
| 3.2.4.3. Giao diện chương trình chính:..... | 57 |
| 3.2.5 Chạy chương trình :..... | 57 |
| 3.2.6 Giao diện quản lý người dùng :..... | 62 |
| 3.2.7 Giao diện quản lý Khoa/Viện:..... | 64 |
| 3.2.8 Giao diện quản lý giảng viên : | 64 |
| 3.2.9 Giao diện quản lý lương :..... | 65 |
| KẾT LUẬN | 66 |

MỞ ĐẦU

Khám phá tri thức - Khai phá dữ liệu (Knowledge discovery - Data mining) là một lĩnh vực quan trọng của ngành Công nghệ thông tin, đã và đang thu hút sự quan tâm đông đảo các nhà khoa học trên thế giới và trong nước tham gia nghiên cứu. Khai phá dữ liệu ra đời vào những năm cuối thập kỷ 80 của thế kỷ XX, nó là lĩnh vực được nghiên cứu nhằm tự động khai thác thông tin, tri thức mới hữu ích, tiềm ẩn từ các CSDL lớn, kho dữ liệu,... Những vấn đề được quan tâm trong khai phá dữ liệu là phân lớp nhận dạng mẫu, luật kết hợp, phân cụm dữ liệu, ... Trong đó, phân cụm dữ liệu (Data Clustering) là một trong những kỹ thuật khai thác dữ liệu có hiệu quả. Phân cụm dữ liệu là quá trình tìm kiếm và phát hiện ra các cụm hoặc các mẫu dữ liệu tự nhiên trong cơ sở dữ liệu lớn. Phân cụm dữ liệu đã được ứng dụng trong nhiều lĩnh vực khác nhau như giáo dục, y tế, kinh tế, bảo hiểm, phân đoạn ảnh, ...

Việc áp dụng phân cụm dữ liệu để phân tích trong ngành kế toán hiện nay là rất cần thiết, bởi lượng dữ liệu lưu trữ lương khá lớn, việc phân tích đánh giá lương để đưa ra các chiến lược cân đối nguồn chi phí của đơn vị, dự báo quỹ lương và có kế hoạch cân đối tài chính cho phù hợp cũng gặp nhiều khó khăn. Ngoài ra việc phân tích lương còn phục vụ công tác quản lý nhân sự, giúp nắm được tình hình sử dụng con người của đơn vị từ đó đưa ra các chính sách tuyển dụng phù hợp, có các giải pháp tạo động lực cho người lao động bằng các chính sách tài chính.

Việc phân cụm dữ liệu để phân tích lương cho kết quả thu được sẽ phân loại theo giá trị lương của mỗi cán bộ, phân loại ra các mức thu nhập cao thấp khác nhau từ đó đưa ra các chính sách cân đối thu chi để có những chính sách ưu đãi phù hợp mà vẫn đảm bảo tài chính của đơn vị.

Với các lý do như vậy tôi chọn đề tài: **“Một số phương pháp phân cụm dữ liệu và ứng dụng trong phân tích lương của cán bộ trường Cao đẳng Nghề Hà Nam”** làm đề tài luận văn tốt nghiệp. Bộ cục luận văn gồm có 3 chương:

Chương I: Tổng quan về khai phá dữ liệu và phân cụm dữ liệu.

Chương II: Một số thuật toán phân cụm dữ liệu điển hình

Chương III: Ứng dụng phương pháp phân nhóm dữ liệu vào phân tích lương của cán bộ trường Cao đẳng Nghề Hà Nam.