

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



NGUYỄN THỊ THÙY DƯƠNG

**NGHIÊN CỨU LÝ THUYẾT NAIVE BAYES VÀ
ỨNG DỤNG TRONG PHÂN LOẠI VĂN BẢN TIẾNG VIỆT**

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 THỊ THU HÀ

THÁI NGUYÊN, NĂM 2015

LỜI CAM ĐOAN

Tôi xin cam đoan:

1. Những nội dung trong luận văn này là do tôi thực hiện dưới sự trực tiếp hướng dẫn của cô giáo TS. Nguyễn Thị Thu Hà.
2. Mọi tham khảo dùng trong luận văn đều được trích dẫn rõ ràng tên tác giả, tên công trình, thời gian, địa điểm công bố.
3. Mọi sao chép không hợp lệ, vi phạm quy chế đào tạo, hay gian trá, tôi xin chịu hoàn toàn trách nhiệm.

Tác giả luận văn

Nguyễn Thị Thùy Dương

LỜI CẢM ƠN

Lời đầu tiên tôi xin được bày tỏ lòng biết ơn chân thành đến Ban Giám Hiệu, các thầy giáo, cô giáo phòng Sau đại học trường Đại học Công Nghệ Thông Tin & Truyền Thông, các thầy giáo ở Viện Công Nghệ Thông Tin đã giảng dạy và tạo mọi điều kiện cho tôi học tập, nghiên cứu và hoàn thành luận văn này.

Đặc biệt, tôi xin bày tỏ sự kính trọng và lòng biết ơn sâu sắc đến TS. Nguyễn Thị Thu Hà, người đã tận tình hướng dẫn và giúp đỡ tôi trong suốt quá trình học tập, nghiên cứu và hoàn thành luận văn.

Tôi chân thành cảm ơn các thầy cô Khoa Công nghệ thông tin, Trường Trung cấp nghề Phát Thanh Truyền Hình Thanh Hóa nơi tôi công tác đã tạo điều kiện và hỗ trợ tôi trong suốt thời gian qua.

Tôi cũng xin chân thành cảm ơn người thân, bạn bè đã giúp đỡ và động viên tôi trong suốt thời gian học tập cũng như trong thời gian thực hiện luận văn.

Xin chân thành cảm ơn!

Thái Nguyên, ngày 20 tháng 08 năm 2015

MỤC LỤC

LỜI CAM ĐOAN	i
LỜI CẢM ƠN	iii
DANH SÁCH CÁC BẢNG.....	vi
Chương 1: TỔNG QUAN VỀ PHÂN LOẠI VĂN BẢN.....	3
1.1. Giới thiệu bài toán phân loại văn bản tiếng Việt	3
1.1.1. Tổng quan bài toán phân loại văn bản	3
1.1.2. Mô hình hệ thống phân loại văn bản.....	4
1.1.3. Các khái niệm cơ bản trong phân loại văn bản.....	5
1.2. Các nghiên cứu liên quan.....	9
1.2.1. Đánh giá phân loại văn bản.....	11
1.2.2. Lý thuyết Naive Bayes.....	11
1.2.3. Khái niệm.....	12
1.3. Kết luận chương 1	17
Chương 2: PHÂN LOẠI VĂN BẢN TIẾNG VIỆT DỰA TRÊN PHƯƠNG PHÁP NAIVE BAYES	18
2.1. Bộ phân loại Naive Bayes.....	18
2.2. Phân loại văn bản tiếng Việt	22
2.2.1. Ứng dụng Naive Bayes trong phân loại văn bản tiếng Việt	22
2.2.2. Rút trích đặc trưng	25
2.2.3. Phân loại văn bản tiếng Việt dựa trên Naive Bayes.....	39
2.3. Kết luận chương 2.....	42
Chương 3: PHÁT TRIỂN HỆ THỐNG PHÂN LOẠI VĂN BẢN TIẾNG VIỆT DỰA TRÊN NAIVE BAYES.....	43
3.1. Mô hình tổng quát của hệ thống	43
3.2. Xây dựng tập ngữ liệu.....	44
3.2.1. Xây dựng tập dữ liệu.....	44
3.2.2. Tiền xử lý và chuẩn hóa dữ liệu.....	47
3.2.3. Xây dựng bộ từ điển danh từ	48
3.3. Môi trường cài đặt	50
3.3.1. Môi trường cài đặt của hệ thống.....	50

3.3.2. Cấu trúc của chương trình.....	50
3.3.3. Giao diện chương trình	51
3.4. Kết quả thực nghiệm.....	56
3.5. Kết luận chương 3.....	57
KẾT LUẬN VÀ HƯỚNG PHÁT TRIỂN	59
TÀI LIỆU THAM KHẢO.....	60

DANH SÁCH CÁC BẢNG

Bảng 1.2. Đánh giá phân loại văn bản	11
Bảng 2.1. Các từ chủ đề trong tập mô tả của Andrews năm 2009.....	30
Bảng 2.2. Danh sách một số chủ đề đã được xây dựng	41
Bảng 3.1. Các chức năng của chương trình	45
Bảng 3.2. Danh sách một số từ trong tập từ chủ đề.	49
Bảng 3.3. Độ triệu hồi khi thực hiện các truy vấn.	57

DANH SÁCH HÌNH VẼ

Hình 1.1. Quá trình học phân loại văn bản.	4
Hình 1.2. Mô hình SVM	8
Hình 2.1. Mô tả bước xây dựng bộ phân lớp	21
Hình 2.2. Trực quan hóa dữ liệu giảm chiều	26
Hình 2.3. Danh sách một số từ dừng.....	27
Hình 2.4. Chỉ số ngữ nghĩa ẩn	28
Hình 2.5. Mô tả việc sắp xếp một văn bản vào chủ đề phù hợp.....	29
Hình 2.6. Mô tả một cách suy diễn chủ đề dựa trên các thuật ngữ.....	30
Hình 2.7. Mô hình chủ đề dựa trên mạng Bayesian	33
Hình 2.8. Mô hình chủ đề dựa trên HMM	34
Hình 2.9. Quy trình phân loại văn bản tiếng Việt.....	36
Hình 2.10. Mô hình chủ đề dựa trên xác suất	37
Hình 2.11. Thuật toán xây dựng mô hình chủ đề.	39
Hình 3.1. Sơ đồ chức năng hệ thống xử lý văn bản tiếng Việt.....	43
Hình 3.2. Biểu đồ Use case tổng quát.....	44
Hình 3.2 Văn bản đã chuẩn hóa.	48
Hình 3.3. Hệ thống VLSP.	49
Hình 3.4. Giao diện trang chủ.....	51
Hình 3.5. Giao diện các thể loại tin	52
Hình 3.6. Giao diện tin huấn luyện	52
Hình 3.7. Giao diện danh sách từ khóa	53
Hình 3.8. Giao diện cài đặt huấn luyện.....	54
Hình 3.9. Giao diện huấn luyện phân loại	55

Hình 3.10. Giao diện danh sách tin tức..... 55

Hình 3.11. Giao diện người dùng 56

DANH SÁCH CÁC CHỮ VIẾT TẮT

Viết tắt	Tiếng Anh	Tiếng Việt
k- NN	k- Nearest Neighbor	k-Láng giềng gần nhất
SVM	Support Vector Machine	Máy véc tơ hỗ trợ
RSS	Really Simple Syndication	Định dạng tập tin
ML	Machine Languages	Ngôn ngữ máy
LSI	Latent Sematic Indexing	Chỉ số ngữ nghĩa ẩn
SVD	Singular Value Decomposition	Phân tích giá trị đơn

MỞ ĐẦU

1. Lý do chọn đề tài

Với lượng thông tin đồ sộ, một yêu cầu lớn đặt ra đối với chúng ta là làm sao tổ chức và tìm kiếm thông tin có hiệu quả nhất. Phân loại thông tin là một trong những giải pháp hợp lý cho yêu cầu trên. Nhưng một thực tế là khối lượng thông tin quá lớn, việc phân loại dữ liệu thủ công là điều không tưởng. Hướng giải quyết là một chương trình máy tính tự động phân loại các thông tin trên.

Đề tài “*Nghiên cứu lý thuyết Naive Bayes và ứng dụng trong phân loại văn bản Tiếng Việt*” nhằm tìm hiểu và thử nghiệm các phương pháp phân loại văn bản áp dụng trên tiếng Việt. Phân loại văn bản (Text classification) là một trong những công cụ khai phá dữ liệu dạng văn bản một cách hữu hiệu, làm nhiệm vụ đưa những văn bản có cùng nội dung chủ đề giống nhau về cùng một lớp có sẵn. Phân loại văn bản giúp người dùng dễ dàng hơn trong việc tìm kiếm thông tin cần thiết đồng thời có thể lưu trữ các thông tin theo đúng chủ đề (topic) hay lớp (class) dựa trên các thuật toán phân loại.

2. Đối tượng và phạm vi nghiên cứu: Tìm hiểu lý thuyết Naive Bayes và ứng dụng trong phân loại văn bản tiếng Việt.

3. Những nội dung nghiên cứu chính

- Chương 1: Tổng quan về phân loại văn bản

Tổng quan về phân loại văn bản và khái niệm cơ bản về lý thuyết Naive Bayes, bộ phân loại Naive Bayes trên mô hình xác suất.

- Chương 2: Phân loại văn bản tiếng Việt dựa trên phương pháp Naive Bayes

Trình bày phương pháp phân loại văn bản tiếng Việt dựa trên phân loại Naive Bayes và cách giảm chiều đặc trưng nhằm tăng tốc trong quá trình tính toán xử lý bằng cách sử dụng mô hình chủ đề dùng cho tiếng Việt.