

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

ĐẶNG ĐÌNH TUYẾN

**PHÂN LỚP VĂN BẢN NHỜ MÁY VÉC – TỜ HỖ TRỢ VỚI HÀM STRING
KERNEL**

Chuyên ngành: Khoa học máy tính

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

Người hướng dẫn khoa học: PGS.TS.Nguyễn Tân Ân

THÁI NGUYÊN - 2016

LỜI CẢM ƠN

Luận văn được hoàn thành tại trường Đại học Công nghệ Thông tin và Truyền thông Thái Nguyên.

Tác giả luận văn xin bày tỏ lòng biết ơn sâu sắc tới thầy hướng dẫn khoa học: PGS.TS Nguyễn Tân Ân đã tận tình hướng dẫn, giúp đỡ và tạo mọi điều kiện để tác giả thực hiện luận văn này. Tác giả cũng xin chân thành cảm ơn tập thể các thầy cô giáo trong khoa CNTT, phòng quản lý sau đại học Trường Đại học Công nghệ Thông tin và Truyền thông Thái Nguyên đã tạo mọi điều kiện giúp đỡ cho tác giả nghiên cứu, học tập và hoàn thành luận văn.

Xin cảm ơn gia đình, bạn bè, đồng nghiệp đã tạo điều kiện thuận lợi về tinh thần và vật chất cho tác giả hoàn thành luận văn này. Xin cảm ơn tất cả!

Thái Nguyên, tháng 6 năm 2016

Tác giả luận văn

Đặng Đình Tuyền

LỜI CAM ĐOAN

Tôi là Đặng Đình Tuyển, học viên cao học K13, chuyên ngành Khoa học máy tính, khoá 2014-2016. Tôi xin cam đoan luận văn thạc sĩ “Phân lớp văn bản nhờ Máy Véc-tơ hỗ trợ (SVM) với hàm string kernel” là công trình nghiên cứu của riêng tôi cùng với sự hướng dẫn của PGS.TS Nguyễn Tân Ân. 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.

Tất cả những tham khảo từ các nghiên cứu liên quan đều được nêu nguồn gốc một cách rõ ràng từ danh mục tài liệu tham khảo của luận văn. Trong luận văn, không có việc sao chép tài liệu, công trình nghiên cứu của người khác mà không chỉ rõ về tài liệu tham khảo.

Thái Nguyên, tháng 6 năm 2016

Tác giả

Đặng Đình Tuyển

MỤC LỤC

LỜI CẢM ƠN	i
LỜI CAM ĐOAN.....	iii
MỤC LỤC.....	iv
DANH MỤC HÌNH ẢNH	vi
DANH MỤC BẢNG BIỂU	vii
CHƯƠNG 1: BÀI TOÁN PHÂN LỚP.....	1
1.1. Nội dung bài toán phân lớp.....	1
1.2. Các phương pháp phân lớp.....	2
1.2.1. Phương pháp Naïve Bayes (NB)	2
1.2.2. Phương pháp K-Nearest Neighbor (kNN)	3
1.2.3. Neural Network (NNet)	5
1.2.4. Centroid- based vector	6
1.3. Máy véc-tơ hỗ trợ (Support Vector Machine SVM).....	7
1.3.1. Bài toán phân loại SVM	7
1.3.2. Ý tưởng của SVM.....	8
1.3.3. Phương pháp tìm α^* , b.	16
1.3.4. SVM đối với bài toán nhiều lớp	21
1.3. Kết luận	24
CHƯƠNG 2: NHỮNG KIẾN THỨC CƠ SỞ	25
2.1. Hàm Kernel	25
2.1.1. Không gian góc, không gian đặc trưng.....	25
2.1.2. Định nghĩa kernel	26
2.1.3. Một số ví dụ về Φ và $k(\cdot)$	26
2.1.4. Một số hàm kernel	28
2.1.5. Định lý	30
2.1.6. Kernel là độ đo giống nhau giữa hai đối tượng	31
2.1.7. Kernel trick	32
2.1.8. Xây dựng các kernel	32

2.1.9. Nhân hóa một số phương pháp phân lớp	34
2.2. String kernel	39
2.2.1. Kernel dựa trên mô hình k_gram	39
2.2.2. Kernel dựa trên trọng số của các xâu	41
2.2.3. Tính string kernel dùng quy hoạch động	43
2.2.4. Kernel dựa trên độ giống nhau giữa hai xâu	44
2.2.5. Một số đặc trưng của Tiếng Việt.	45
2.3. Kết luận	48
CHƯƠNG 3: CÀI ĐẶT THỬ NGHIỆM THUẬT TOÁN SVM CHO BÀI TOÁN	
TÌM KIẾM VĂN BẢN	49
3.1. Mô tả bài toán	49
3.2. Phân tích, cài đặt thuật toán	49
3.2.1. Thuật toán huấn luyện để tìm từ khóa	49
3.2.2. Thuật toán sử dụng từ khóa tìm kiếm văn bản	57
3.3. Kết luận	61
KẾT LUẬN VÀ HƯỚNG PHÁT TRIỂN	62
TÀI LIỆU THAM KHẢO	63

DANH MỤC HÌNH ẢNH

Hình 1.1: Kiến trúc mô đun (Modular Architecture). Các kết quả của từng mạng con sẽ là giá trị đầu vào cho mạng siêu chủ đề và được nhân lại với nhau để dự đoán chủ đề cuối cùng.....	6
Hình 1.2: Các trường hợp của siêu mặt h phân chia tập dữ liệu D trong SVM.....	8
Hình 1.3: Siêu mặt phân chia tập mẫu huấn luyện với 2 lớp là lớp + 1 hình vuông và lớp – 1 hình tròn.....	9
Hình 1.4: Siêu phẳng tuyến tính phân chia dữ liệu, m là khoảng cách giữa hai lề...	10
Hình 1.5: Nguyên lý cơ bản của phương pháp một-chơi-phần còn lại cho ba lớp ...	22
Hình 1.6: Nguyên lý cơ bản của phương pháp phân chia một-chơi-một.....	22
Hình 1.7: Biểu diễn phương pháp END để phân chia ba trạng thái của bài toán dự đoán trong phân lớp.....	24
Hình 2.1: Mỗi điểm dữ liệu được ánh xạ bằng một hàm không tuyến tính Φ từ không gian dữ liệu X vào không gian đặc trưng F. Trong đó $\Phi(x)$ và $\Phi(o)$ là các véc-tơ đặc trưng của các điểm dữ liệu gốc x và o.....	26
Hình 2.2: Ánh xạ dữ liệu từ không gian đầu vào R^2 sang không gian dữ liệu R^3	27
Hình 2.3: Kernel đa thức bậc hai ánh xạ từ không gian hai chiều vào không gian đặc trưng 3 chiều.....	29
Hình 2.4: Dữ liệu được tách thành hai lớp trong không gian ban đầu.....	31
Hình 3.1: Trang web Du lịch Khát vọng Việt.....	50
Hình 3.2: Trang web taxinoibaiphuonglong.com	52
Hình 3.3: Trang web vietnamtourism.com	55

DANH MỤC BẢNG BIỂU

Bảng 3.1: Bảng thống kê các từ đặc trưng từ Đoạn mẫu 1	50
Bảng 3.2: Tính toán tần xuất và trọng số của các từ (theo định nghĩa từ tiếng Việt).....	51
Bảng 3.3: Bảng thống kê các từ đặc trưng từ Đoạn mẫu 2.	52
Bảng 3.4: Tính toán tần xuất và trọng số của các từ (theo định nghĩa từ tiếng Việt).....	54
Bảng 3.5: Bảng thống kê các từ đặc trưng từ Đoạn mẫu 3.	55
Bảng 3.6: Tính toán tần xuất và trọng số của các từ (theo định nghĩa từ tiếng Việt).....	56
Bảng 3.7: Bảng tổng hợp.....	56
Bảng 3.8: Số lần xuất hiện của các từ trong các văn bản huấn luyện	59
Bảng 3.9: Bảng phân nhóm với nhãn là “Vịnh Hạ Long”	59
Bảng 3.10: Bảng phân nhóm với nhãn là “Di sản”	60
Bảng 3.11: Bảng phân nhóm với nhãn là “đảo”	60

CHƯƠNG 1: BÀI TOÁN PHÂN LỚP

1.1. Nội dung bài toán phân lớp

Phân lớp (classification) là một tiến trình xử lý nhằm xếp các mẫu dữ liệu hay các đối tượng vào một trong các lớp đã được định nghĩa trước. Các mẫu dữ liệu hay các đối tượng được xếp về các lớp dựa vào giá trị của các thuộc tính (attributes) cho một mẫu dữ liệu hay đối tượng. Sau khi đã xếp tất cả các đối tượng đã biết trước vào các lớp tương ứng, lúc này mỗi lớp được đặc trưng bởi tập các thuộc tính của các đối tượng chứa trong lớp đó. Ví dụ: phân lớp văn bản, tế bào để xác định tế bào ung thư.

Phân lớp còn được gọi là phân lớp có giám sát (supervised classification), là một trong những lĩnh vực phổ biến nhất của học máy (machine learning) và khai thác dữ liệu (data mining). Nó giải quyết việc xác định những quy tắc giữa số lượng biến số độc lập và kết quả đạt được hay một biến số xác định phụ thuộc trong tập dữ liệu được đưa ra. Tổng quát, đưa ra một tập mẫu học $(x_1, x_2, \dots, x_k, y_i)$, $i=1, \dots, N$, nhiệm vụ là phải ước lượng được một bộ phân lớp hay một mô hình xấp xỉ một hàm $y = f(x)$ chưa biết mà phân lớp chính xác cho bất kỳ mẫu nào thuộc tập các mẫu học. Có nhiều cách để biểu diễn một mô hình phân lớp và có rất nhiều thuật toán giải quyết nó. Các thuật toán phân lớp tiêu biểu bao gồm như mạng neural, cây quyết định, suy luận quy nạp, mạng Bayesian, Support Vector Machine.... Tất cả các cách tiếp cận này xây dựng những mô hình đều có khả năng phân lớp cho một mẫu mới chưa biết dựa vào những mẫu tương tự đã được học.

Bài toán phân lớp có thể xử lý thông tin được thu thập từ mọi lĩnh vực hoạt động của con người và thế giới tự nhiên được biểu diễn dưới dạng các bảng. Bảng này bao gồm các đối tượng và các thuộc tính. Các phần tử trong bảng là các giá trị xác định các thuộc tính (attributes hay features) của các đối tượng. Trong đó số cột chính là số thuộc tính của các đối tượng, mỗi cột là một thuộc tính và số dòng chính là số đối tượng chứa trong dữ liệu này. Mọi dữ liệu được biểu diễn dưới các dạng khác có thể được chuyển thành dạng bảng như trên để thực hiện quá trình phân lớp.

1.2. Các phương pháp phân lớp

1.2.1. Phương pháp Naïve Bayes (NB)

NB là phương pháp phân loại dựa vào xác suất được sử dụng rộng rãi trong lĩnh vực máy học (Mitchell trình bày năm 1996, Joachims trình bày năm 1997 và Jason năm 2001) được sử dụng lần đầu tiên trong lĩnh vực phân loại bởi Maron vào năm 1961, sau đó trở nên phổ biến dùng trong nhiều lĩnh vực như trong các công cụ tìm kiếm (được mô tả năm 1970 bởi Rijsbergen), các bộ lọc mail (mô tả năm 1998 bởi Sahami)...

* Ý tưởng

Ý tưởng cơ bản của cách tiếp cận Naïve Bayes là sử dụng xác suất có điều kiện giữa từ và chủ đề để dự đoán xác suất chủ đề của một văn bản cần phân loại. Điểm quan trọng của phương pháp này chính là ở chỗ giả định rằng sự xuất hiện của tất cả các từ trong văn bản đều độc lập với nhau. Với giả định này NB không sử dụng sự phụ thuộc của nhiều từ vào một chủ đề, không sử dụng việc kết hợp các từ để đưa ra phán đoán chủ đề và do đó việc tính toán NB chạy nhanh hơn các phương pháp khác với độ phức tạp theo hàm số mũ.

* Công thức

Mục đích chính là tính được xác suất $Pr(C_j, d')$, xác suất để văn bản d' nằm trong lớp C_j . Theo luật Bayes, văn bản d' sẽ được gán vào lớp C_j nào có xác suất $Pr(C_j, d')$ cao nhất. Công thức sau dùng để tính $Pr(C_j, d')$ (do Joachims đề xuất năm 1997)

$$H_{BAYES} = \arg \max_{C_j \in C} \left(\frac{\Pr(C_j) \cdot \prod_{i=1}^{d'} \Pr(w_i | C_j)}{\sum_{C' \in C} \Pr(C') \cdot \prod_{i=1}^{d'} \Pr(w_i | C')} \right) = \arg \max_{C_j \in C} \left(\frac{\Pr(C_j) \cdot \prod_{i=1}^{d'} \Pr(w_i | C_j)^{IF(w_i, d')}}{\sum_{C' \in C} \Pr(C') \cdot \prod_{i=1}^{d'} \Pr(w_i | C')^{IF(w_i, d')}} \right)$$

Với:

- (TF, d') là số lần xuất hiện của từ w_i trong văn bản d'
- $|d'|$ là số lượng các từ trong văn bản d'
- w_i là một từ trong không gian đặc trưng F với số chiều là $|F|$

Số hóa bởi Trung tâm Học liệu – ĐHTN <http://www.lrc.tnu.edu.vn>

$\Pr(C_j)$ được tính dựa trên tỷ lệ phần trăm của số văn bản mỗi lớp

$$\Pr(C_j) = \frac{\|C_j\|}{\|C\|} = \frac{\|C_j\|}{\sum_{C' \in C} \|C'\|}$$

tương ứng với tập dữ liệu huấn luyện.

- $\Pr(w_i|C_j)$ được tính sử dụng phép ước lượng Laplace (do Naplik trình bày năm 1982)

$$\Pr(w_i | C_j) = \frac{1 + TF(w_i, C_j)}{|F| + \sum_{w' \in |F|} TF(w', C_j)}$$

Ngoài ra còn có các phương pháp NB khác có thể kể ra như sau ML Naive Bayes, MAP Naive Bayes, Expected Naive Bayes, Bayesian Naive Bayes (Jason mô tả năm 2001). Naive Bayes là một công cụ rất hiệu quả trong một số trường hợp. Kết quả có thể rất tồi nếu dữ liệu huấn luyện nghèo nàn và các tham số dự đoán (như không gian đặc trưng) có chất lượng kém. Nhìn chung đây là một thuật toán phân loại tuyến tính thích hợp trong phân loại văn bản nhiều chủ đề. NB có ưu điểm là cài đặt đơn giản, tốc độ nhanh, dễ dàng cập nhật dữ liệu huấn luyện mới và có tính độc lập cao với tập huấn luyện, có thể sử dụng kết hợp nhiều tập huấn luyện khác nhau. Tuy nhiên NB ngoài giả định tính độc lập giữa các từ còn phải cần đến một ngưỡng tối ưu để cho kết quả khả quan. Nhằm mục đích cải thiện hiệu năng của NB, các phương pháp như multiclass-boosting, ECOC (do Berger trình bày năm 1999 và Ghani mô tả lại năm 2000) có thể được dùng kết hợp.

1.2.2. Phương pháp K-Nearest Neighbor (kNN)

Đây là phương pháp truyền thống khá nổi tiếng về hướng tiếp cận dựa trên thống kê đã được nghiên cứu trong nhận dạng mẫu hơn bốn thập kỷ qua (theo tài liệu của Dasarathy năm 1991). kNN được đánh giá là một trong những phương pháp tốt nhất (áp dụng trên tập dữ liệu Reuters phiên bản 21450), được sử dụng từ những thời kỳ đầu của việc phân loại văn bản (được trình bày bởi Marsand năm 1992, Yang năm 1994, Iwayama năm 1995)