

**HỌC VIỆN CÔNG NGHỆ BƯU CHÍNH VIỄN THÔNG**

-----



**LÊ TRUNG HIẾU**

**DỰ ĐOÁN GIỚI TÍNH NGƯỜI DÙNG INTERNET DỰA TRÊN  
LỊCH SỬ TRUY CẬP**

**Chuyên ngành:     Hệ thống thông tin**

**Mã số:               60.48.01.04**

**TÓM TẮT LUẬN VĂN THẠC SĨ**

**HÀ NỘI - 2017**

Luận văn được hoàn thành tại:  
**HỌC VIỆN CÔNG NGHỆ BƯU CHÍNH VIỄN THÔNG**

Người hướng dẫn khoa học: .....PGS.TS Từ Minh Phương.....

Phản biện 1: .....

Phản biện 2: .....

Luận văn sẽ được bảo vệ trước Hội đồng chấm luận văn thạc sĩ tại Học viện  
Công nghệ Bưu chính Viễn thông

Vào lúc: ..... giờ ..... ngày ..... tháng ..... .. năm .....

Có thể tìm hiểu luận văn tại:

- Thư viện của Học viện Công nghệ Bưu chính Viễn thông

## MỤC LỤC

<b>MỤC LỤC .....</b>	<b>1</b>
<b>MỞ ĐẦU .....</b>	<b>3</b>
<b>CHƯƠNG 1: TỔNG QUAN VỀ DỰ ĐOÁN GIỚI TÍNH NGƯỜI DÙNG INTERNET .....</b>	<b>5</b>
<b>1.1 Bài toán xác định giới tính và ứng dụng của bài toán vào thực tiễn.....</b>	<b>5</b>
1.1.1 Mở đầu.....	5
1.1.2 Bài toán xác định giới tính.....	5
1.1.3 Ứng dụng của bài toán vào thực tiễn.....	6
<b>1.2 Các dạng dữ liệu lịch sử có thể dự đoán.....</b>	<b>7</b>
<b>1.3 Các phương pháp xác định giới tính đã có .....</b>	<b>7</b>
1.3.1 Phương pháp xác định giới tính sử dụng bài viết từ blog.....	7
1.3.2 Phương pháp xác định giới tính sử dụng dữ liệu thông tin di động liên lạc hàng ngày .....	8
1.3.3 Xác định giới tính sử dụng dữ liệu từ các thông điệp trên twitter bằng phương pháp hồi quy.....	8
<b>1.4 Kết luận chương .....</b>	<b>9</b>
<b>CHƯƠNG 2: DỰ ĐOÁN GIỚI TÍNH NGƯỜI DÙNG INTERNET SỬ DỤNG LỊCH SỬ TRUY CẬP .....</b>	<b>10</b>
<b>2.1 Giới thiệu về phương pháp học máy SVM .....</b>	<b>10</b>
2.1.1 Giới thiệu về SVM .....	10
2.1.2 Bài toán phân 2 lớp với SVM.....	11
2.1.3 Các bước chính của phương pháp SVM.....	11
2.1.4 Ưu điểm phương pháp SVM trong phân lớp dữ liệu .....	12
<b>2.2 Giới thiệu về dữ liệu sử dụng .....</b>	<b>12</b>
<b>2.3 Các dạng đặc trưng dùng trong phân lớp.....</b>	<b>13</b>
2.3.1 Dạng đặc trưng theo mốc thời gian.....	13
2.3.2 Dạng đặc trưng về danh mục và chủng loại sản phẩm .....	13
<b>2.4 Xây dựng mô hình dự đoán giới tính dựa trên học máy có giám sát .....</b>	<b>14</b>
2.4.1 Tiền xử lý dữ liệu .....	15
2.4.2 Biểu diễn dữ liệu .....	15
<b>2.5 Kết luận chương .....</b>	<b>16</b>
<b>CHƯƠNG 3: THỬ NGHIỆM VÀ ĐÁNH GIÁ .....</b>	<b>17</b>
<b>3.1 Mô tả dữ liệu.....</b>	<b>17</b>

<b>3.2 Các tiêu chuẩn đánh giá .....</b>	<b>17</b>
<b>3.3 Phương pháp thực nghiệm.....</b>	<b>17</b>
<b>3.4 Kết quả thực nghiệm.....</b>	<b>19</b>
<b>3.5 So sánh với một số phương pháp khác .....</b>	<b>20</b>
<b>3.6 Kết luận chương.....</b>	<b>22</b>
<b>KẾT LUẬN .....</b>	<b>23</b>
<b>TÀI LIỆU THAM KHẢO .....</b>	<b>24</b>

## MỞ ĐẦU

Ngày nay, người ta thường dành một lượng lớn thời gian trong ngày để truy cập Internet. Internet được người dùng sử dụng cho việc tìm kiếm thông tin, đọc tin tức, mua sắm, chơi trò chơi v.v. Và các nhà quảng cáo không thể bỏ lỡ cơ hội để tiếp thị trực tuyến đến với khách hàng của họ nhằm cung cấp các dịch vụ phù hợp với nhu cầu của tổ chức, cá nhân sử dụng mạng Internet. Tuy nhiên, hiện nay các nhà quảng cáo đang cung cấp toàn bộ thông tin của mình đến tất cả khách hàng họ có. Chính vì vậy người dùng thường phải đối mặt với số lượng lớn các thông tin không phù hợp ví dụ như không phù hợp về độ tuổi, về nghề nghiệp, về văn hóa và giới tính.

Tình trạng quá tải thông tin không đến đích này dẫn đến sự sụt giảm đáng kể trong việc tiếp thị trực tuyến. Từ đó việc phân loại người dùng Internet để đưa ra các số liệu thống kê, kế hoạch quảng cáo giúp hệ thống tiếp cận cung cấp thông tin phù hợp, hữu ích cho từng đối tượng tương đối quan trọng. Xuất phát từ thực trạng đang xảy ra, luận văn sẽ trình bày về phương pháp xác định giới tính để phân loại người dùng Internet được thực hiện bằng kỹ thuật học máy, sử dụng thông tin người dùng đã biết giới tính và các thông tin về lịch sử truy cập web của họ để huấn luyện máy nhận biết giới tính của những người dùng khác khi ta chỉ biết lịch sử truy cập các trang web và dữ liệu danh mục mà người đó quan tâm.

Với mục tiêu đặt ra như vậy, nội dung và kết quả của luận văn được trình bày qua 3 chương như sau:

Chương 1 giới thiệu về dữ liệu truy cập của người dùng Internet thông qua thống kê, các khái niệm và đặc trưng trong tập dữ liệu này, bao gồm các mối quan hệ giữa các trang thông tin và người dùng mạng, những hành vi của người dùng khi truy cập Internet, cách thức truy cập, tìm kiếm thông tin. Giới thiệu những phương pháp nhằm mục tiêu theo hành vi hiện nay được áp dụng cho người dùng Internet và những hạn chế của các phương pháp này.

Chương 2 trình bày tổng quan về kỹ thuật học máy, một số kỹ thuật học máy và tập trung vào kỹ thuật được sử dụng trong luận văn là kỹ thuật học máy SVM. Dựa vào những đặc trưng việc truy cập thông tin của người dùng Internet, đưa ra phương pháp dự đoán giới tính áp dụng kỹ thuật học máy và xếp hạng tỉ lệ độ chính xác nhằm tăng hiệu quả dự đoán so với các phương pháp đang tồn tại.

Chương 3 trình bày kết quả thực nghiệm và đánh giá. Sử dụng dữ liệu có sẵn PAKDD'15 được cung cấp bởi Công ty Cổ phần FPT (<http://www.fpt.com.vn>), thực hiện xây dựng bộ dữ liệu từ dữ liệu thực tế chưa chuẩn hóa hiện có PAKDD'15 cho một số lượng người dùng, sử dụng kỹ thuật học máy SVM ở chương 2 và một số công cụ để đưa ra tỉ lệ, độ chính xác của phương pháp dự đoán giới tính dựa trên lịch sử truy cập. Đánh giá kết quả so với các phương pháp dự đoán khác, và so sánh với cách làm việc hiện tại trong việc dự đoán giới tính.

## **CHƯƠNG 1: TỔNG QUAN VỀ DỰ ĐOÁN GIỚI TÍNH NGƯỜI DÙNG INTERNET**

### **1.1 Bài toán xác định giới tính và ứng dụng của bài toán vào thực tiễn**

#### **1.1.1 Mở đầu**

Ngày nay, với sự phát triển không ngừng của khoa học công nghệ trên thế giới nói chung và ở Việt Nam nói riêng có những bước tiến vượt bậc. Cơ sở hạ tầng và các trang thiết bị tương đối hiện đại và không ngừng phát triển. Theo báo cáo tổng kết của Bộ TT&TT năm 2016, tỷ lệ người sử dụng Internet ở Việt Nam đạt 62,76% dân số, trong đó tỷ lệ hộ gia đình có truy cập Internet đạt 24,38%, tức là cứ 5 gia đình thì có một hộ sử dụng băng thông rộng cố định. Trong đó, theo thống kê của Cục Viễn thông (Bộ TT&TT) tháng 11/2016, tổng số thuê bao Internet băng rộng cố định đạt hơn 9 triệu thuê bao và số thuê bao băng rộng di động đạt hơn 12,6 triệu thuê bao.

Bên cạnh đó, theo thống kê của “wearesocial.net”, tháng 1-2015, người Việt Nam đang đứng thứ 4 trên thế giới về thời gian sử dụng Internet với 5,2 giờ mỗi ngày, chỉ sau Philippines đứng đầu là 6 giờ, tiếp đó là Thái Lan với 5,5 giờ, và Brazil là 5,4 giờ/ngày.

Chính vì sự phát triển không ngừng của công nghệ thông tin và mức độ phổ biến của Internet ngày nay mà thông tin đến với người dùng vô cùng phong phú và liên tục. Người sử dụng Internet hiện nay thường có thói quen truy cập và tìm kiếm đến những các vấn đề mình quan tâm. Hầu hết các thông tin được lưu vào như một phiên làm việc trên mạng. Các thông tin đó có thể là các bài báo, các tài liệu kinh doanh, sản phẩm, các thông tin kinh tế, thương mại điện tử, các thông tin cá nhân khác. Từ đó đã xuất hiện các nhu cầu phân tích thông tin để phân loại các thông tin đó cho các mục đích khác nhau như học tập, nghiên cứu, kinh doanh, tiếp thị thương mại... Với thực tế đó, ta phải xác định và phân loại những thông tin hữu ích từ các nguồn dữ liệu phong phú từ các phiên làm việc, sử dụng Internet của người dùng sao cho phù hợp với đối tượng cụ thể và hỗ trợ các công cụ tự động hoá trợ giúp trong việc phát hiện tri thức và khai thác thông tin.

#### **1.1.2 Bài toán xác định giới tính**

Việc sử dụng các hoạt động và truy cập Internet có sự khác nhau giữa nam giới và nữ giới. Trung bình một ngày nam giới dành thời gian nhiều hơn cho Internet. Nam giới cũng có một số hoạt động trực tuyến giống với nữ giới. Tuy nhiên có những khác nhau cụ thể ví dụ như nam giới có khuynh hướng truy cập những đặc trưng như tin tức thời sự, bóng đá, hay

trò chơi và các mặt hàng dành cho nam giới. Trái lại nữ giới thường thích thú với các mục mua sắm, thương mại điện tử, chat và tham gia các trang mạng xã hội và blog.

Dự đoán giới tính (hay Determination Gender hoặc Gender Prediction) là phương pháp phân loại và xác định các hoạt động được truy cập bởi giới tính Nam hoặc giới tính Nữ từ những hoạt động khác đã biết trước nhãn. Ví dụ một bài báo trong một trang web có thể được truy cập bởi giới tính nam hoặc giới tính nữ (như thể thao, giáo dục, pháp luật, công nghệ thông tin, mỹ phẩm, quần áo ...). Việc phân loại có thể được tiến hành một cách thủ công: đọc nội dung của từng hoạt động và gán nó vào một nhãn nào đó. Tuy nhiên, đối với hệ thống gồm rất bản ghi thì phương pháp này sẽ tốn rất nhiều thời gian và công sức. Do vậy cần phải có phương pháp tự động để phân loại giới tính. Phương pháp này giúp cho việc xác định giới tính đạt độ chính xác cao và sử dụng cho các mục đích như học tập, nghiên cứu, kinh doanh, tiếp thị thương mại.

Để tiến hành phân loại xác định giới tính nói chung, chúng ta sẽ thực hiện các bước sau đây:

- Bước 1: Xây dựng bộ dữ liệu huấn luyện dựa trên tập dữ liệu thu thập của người dùng đã được phân loại sẵn. Tiến hành học cho bộ dữ liệu, xử lý và thu thập được dữ liệu của quá trình học là các đặc trưng riêng biệt cho từng nội dung.
- Bước 2: Dữ liệu cần phân loại được xử lý, rút ra các đặc trưng kết hợp với đặc trưng được học trước đó để phân loại và đưa ra kết quả.

Đặc điểm nổi bật của bài toán này là sự đa dạng của hoạt động và đặc trưng của nam giới và nữ giới. Các đặc trưng làm cho sự phân loại chỉ mang tính tương đối và có phần chủ quan, nếu do con người thực hiện có thể dễ bị nhập nhằng. Ví dụ có hoạt động truy cập về xem thông tin mua sắm quần áo tại một trang web thương mại điện tử, hoạt động truy cập này vẫn có thể được truy cập bởi nam giới hoặc nữ giới.

### **1.1.3 Ứng dụng của bài toán vào thực tiễn**

Hiện nay, công nghệ ngày càng phát triển, đặc biệt với sự ra đời của các trang mạng xã hội, thương mại điện tử nên lượng thông tin lớn, phi cấu trúc, phức tạp, thậm chí là các thông tin rác cũng rất nhiều. Hầu hết các thông tin đều là các hoạt động trực tuyến như tìm kiếm thông tin, chat, email, mua sắm trực tuyến ... Từ thực tế đó đã xuất hiện các nhu cầu phân tích thông tin của người dùng Internet để phân loại các thông tin đó sao cho phù hợp với



giới tính nhằm đưa ra các số liệu thống kê, kế hoạch quảng cáo giúp hệ thống tiếp cận cung cấp thông tin phù hợp, hữu ích cho từng đối tượng.

## **1.2 Các dạng dữ liệu lịch sử có thể dự đoán**

Có nhiều loại dữ liệu lịch sử có thể được sử dụng để dự đoán. Ở giai đoạn đầu phân loại giới tính, hầu hết các nghiên cứu về lĩnh vực này tập trung vào việc nghiên cứu tác giả, đó là những nhiệm vụ xác định hoặc dự đoán các đặc điểm tác giả bằng cách phân tích các câu chuyện, tác phẩm, tiểu thuyết được tạo ra bởi tác giả nam hay tác giả nữ. Các phương pháp mà các nhà nghiên cứu sử dụng trong các nghiên cứu này chủ yếu dựa trên việc phân tích các phong cách viết, văn phong sử dụng các đặc trưng về ngữ pháp chẳng hạn như từ vựng, cú pháp, hoặc các đặc trưng dựa trên nội dung. Nghiên cứu đầu tiên trong lĩnh vực này bắt đầu vào thế kỷ 19 khi Mendenhall (1887) [16] đã nghiên cứu các tác phẩm của Shakespeare.

Gần đây, do sự phát triển của Internet và các kênh truyền thông trực tuyến, các dạng dữ liệu được thu thập chủ yếu dựa trên nội dung truyền thông ví dụ như: Email, Blog, Twitter, Facebook ...

## **1.3 Các phương pháp xác định giới tính đã có**

### **1.3.1 Phương pháp xác định giới tính sử dụng bài viết từ blog**

Trong những năm trở về trước, Blog là một loại nhật ký, website cá nhân phổ biến chia sẻ những kinh nghiệm sống hoặc một thông tin gì đó trong cuộc sống hằng ngày của con người. Đây là một loại dữ liệu rất rất lớn chứa các bài viết, văn bản do hàng trăm nghìn tác giả người dùng tạo ra. Những thông tin này chứa đựng rất nhiều các đặc trưng có thể khai thác cho bài toán phân loại, cụ thể ở đây là việc xác định giới tính các blogger. Bài báo nghiên cứu cụ thể về xác định nhân khẩu học và giới tính được Schler et al [19] xây dựng năm 2007 với tập dữ liệu là tất cả blog được truy cập trong một ngày tháng 8 năm 2004. Nội dung nghiên cứu chú trọng sự khác biệt trong việc viết blog và sự khác biệt giữa nam giới và nữ giới giữa các blogger ở các độ tuổi khác nhau. Các đặc trưng về phong cách và nội dung được đưa ra làm hạt nhân để giải quyết bài toán.

Các kết quả kiểm thử cho thấy được việc phân loại được các blogger theo giới tính theo các nhóm tuổi, kiểu viết và nội dung. Trong các trường hợp được đưa ra, thì sự kết hợp của các đặc trưng phong cách và nội dung cung cấp độ chính xác phân loại tốt nhất.

### 1.3.2 Phương pháp xác định giới tính sử dụng dữ liệu thông tin di động liên lạc hàng ngày

Phương pháp xác định giới tính thông qua dữ liệu từ các thông tin di động liên lạc hàng ngày được nghiên cứu theo bài báo Demographic Prediction Based on User's Mobile Behaviors [14] trong cuộc thi MDC Data Set. Trong bài báo này, nhóm nghiên cứu đề xuất một mô hình mới cụ thể là Multi-Level Classification Model (Mô hình phân loại Đa cấp) để giải quyết vấn đề các lớp không cân bằng hiện có trong dữ liệu. Dựa trên mô hình này, sẽ đưa ra kết quả việc dự đoán giới tính của người dùng bằng cách kết hợp nhiều mô hình phân loại vào một cấu trúc đa cấp.

Như đã đề cập, tài nguyên dữ liệu hiện có là dữ liệu nhật ký điện thoại di động của người dùng các vị trí khác nhau và thời gian khác nhau. Do đó, nghiên cứu chú trọng các đặc trưng hành vi người dùng và tìm kiếm các đặc trưng độc đáo của các vị trí được ghi lại trong nhật ký di động của tập dữ liệu MDC. Tập dữ liệu được trích xuất phân loại huấn luyện và phân chia theo các tầng, từ tầng 1 đến tầng thấp hơn, lần lượt xác định phân loại ở mỗi tầng cho đến khi thu được kết quả phân loại chính xác nhất.

### 1.3.3 Xác định giới tính sử dụng dữ liệu từ các thông điệp trên twitter bằng phương pháp hồi quy

Xác định giới tính sử dụng dữ liệu từ các thông điệp Twitter là phương pháp phân loại cho từng bình luận theo đặc trưng dựa trên nội dung bình luận bằng phương pháp hồi quy. Ở bước đầu tiên, từ tập dữ liệu thô là những ý kiến trên Twitter được thu thập theo chủ đề, ta tiến hành tiền xử lý các kí tự đặc biệt của Twitter, các kí tự trùng lặp gần nhau, từ viết tắt, tiếng lóng, biểu tượng cảm xúc, mạng ngữ nghĩa. Nghiên cứu được trình bày bởi Dong Nguyen và các cộng sự năm 2013 [13].

Hồi Quy (regression) là một phương pháp học có giám sát (supervised learning) trong Máy Học. Mục tiêu chính là tìm ra mối quan hệ giữa các đặc trưng của một vấn đề nào đó. Cụ thể hơn, từ một tập dữ liệu cho trước, ta xây dựng một mô hình (phương trình, đồ thị, ...) khớp nhất với tập dữ liệu, thể hiện được xu hướng biến thiên và mối quan hệ giữa các đặc trưng. Khi có một mẫu dữ liệu mới vào, dựa vào mô hình, chúng ta có thể dự đoán giá trị của mẫu dữ liệu đó. Lấy ví dụ như chúng ta cần dự đoán giới tính của một twitter dựa vào nội dung và đặc trưng viết của twitter đó. Như vậy chúng ta cần tìm mối quan hệ giữa giới tính phụ thuộc vào nội dung và đặc trưng viết. Dựa vào tập dữ liệu (giả sử thu thập nội dung,

đặc trưng viết và các ký tự đặc biệt của 100 người dùng twitter), ta xây dựng một phương trình  $y = \theta_0 + \theta_1 x_1 + \theta_2 x_2$  trong đó  $y$  là giới tính phụ thuộc  $x_1$  (nội dung) và  $x_2$  (đặc trưng viết). Khi có thêm một mẫu dữ liệu của một người dùng mới, chỉ cần áp vào phương trình như vậy ta sẽ dự đoán được giới tính của người đó.

#### 1.4 Kết luận chương

Chương này đã giới thiệu tổng quan về bài toán xác định giới tính, ứng dụng của bài toán vào thực tiễn và một số phương pháp xác định giới tính và dữ liệu lịch sử liên quan đến việc phân loại giới tính nam hay giới tính nữ. Bên cạnh đó, chương 1 còn đưa ra lý do và thực trạng các hoạt động của người dùng Internet trong luận văn. Ngoài ra cần lưu ý đến yếu tố quan trọng tác động đến kết quả phân loại giới tính đó là phải có một tập dữ liệu lịch sử để huấn luyện chuẩn và đủ lớn để cho thuật toán học phân loại. Nếu chúng ta có được một tập dữ liệu chuẩn và đủ lớn thì quá trình huấn luyện sẽ tốt và khi đó chúng ta sẽ có kết quả phân loại tốt sau khi đã được học. Trong chương 1, luận văn cũng đã giới thiệu một số phương pháp xác định giới tính đã được nghiên cứu trong thời gian gần đây. Những mô tả của chương 1 sẽ làm tiền đề cho việc xác định giới tính người dùng Internet sử dụng dữ liệu lịch sử truy cập trong các chương tiếp theo.

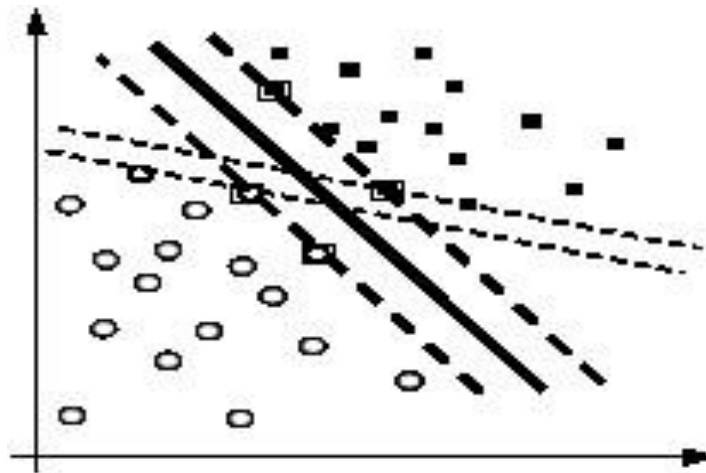
## CHƯƠNG 2: DỰ ĐOÁN GIỚI TÍNH NGƯỜI DÙNG INTERNET SỬ DỤNG LỊCH SỬ TRUY CẬP

### 2.1 Giới thiệu về phương pháp học máy SVM

#### 2.1.1 Giới thiệu về SVM

Cho trước một tập huấn luyện, được biểu diễn trong không gian vector, trong đó mỗi tài liệu là một điểm, phương pháp này tìm ra một siêu phẳng  $f$  quyết định tốt nhất có thể chia các điểm trên không gian này thành hai lớp riêng biệt tương ứng là lớp “+” và lớp “-”. Chất lượng của siêu phẳng này được quyết định bởi khoảng cách (gọi là biên) của điểm dữ liệu gần nhất của mỗi lớp đến mặt phẳng này. Khi đó, khoảng cách biên càng lớn thì mặt phẳng quyết định càng tốt, đồng thời việc phân loại càng chính xác. Ý tưởng của nó là ánh xạ (tuyến tính hoặc phi tuyến) dữ liệu vào không gian các vector đặc trưng (space of feature vectors) mà ở đó một siêu phẳng tối ưu được tìm ra để tách dữ liệu thuộc hai lớp khác nhau.

Mục đích của phương pháp SVM là tìm được khoảng cách biên lớn nhất:



Hình 2.1 Mô tả phương pháp SVM

Đường tô đậm là siêu phẳng tốt nhất và các điểm được bao bởi hình chữ nhật là những điểm gần siêu phẳng nhất, chúng được gọi là các vector hỗ trợ (support vector). Các đường nét đứt mà các support vector nằm trên đó được gọi là lề (margin).

SVM thực chất là một bài toán tối ưu, mục tiêu của thuật toán này là tìm được một không gian  $F$  và siêu phẳng quyết định  $f$  trên  $F$  sao cho sai số phân loại là thấp nhất.

Cho tập mẫu  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_l, y_l)\}$  với  $x_i \in \mathbb{R}^n$ , thuộc vào hai lớp nhãn  $y_i \in \{-1, 1\}$  là nhãn lớp tương ứng của các  $x_i$  (-1 biểu thị lớp I, 1 biểu thị lớp II).

Ta có, phương trình siêu phẳng chứa vector  $\vec{x}_i$  trong không gian:

$$\vec{x}_i \cdot \vec{w} + b = 0$$

Đặt:

$$f(\vec{x}_i) = \text{sign}(\vec{x}_i \cdot \vec{w} + b) = \begin{cases} +1, \vec{x}_i \cdot \vec{w} + b > 0 \\ -1, \vec{x}_i \cdot \vec{w} + b < 0 \end{cases}$$

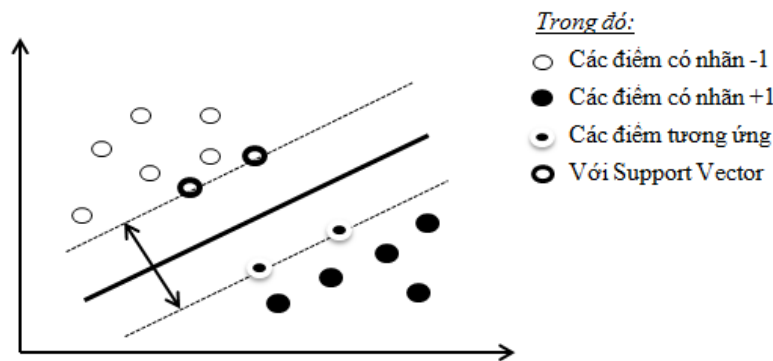
Như vậy,  $f(\vec{x}_i)$  biểu diễn sự phân lớp của  $\vec{x}_i$  vào hai lớp như đã nêu.

Ta nói  $y_i = +1$  nếu  $\vec{x}_i$  thuộc lớp I và  $y_i = -1$  nếu  $\vec{x}_i$  thuộc lớp II.

### 2.1.2 Bài toán phân 2 lớp với SVM

Bài toán đặt ra là: Xác định hàm phân lớp để phân lớp các mẫu trong tương lai, nghĩa là với một mẫu dữ liệu mới  $x_i$  thì cần phải xác định  $x_i$  được phân vào lớp  $+1$  hay lớp  $-1$ .

Tập D có thể phân chia tuyến tính được mà không có nhiễu (tất cả các điểm được gán nhãn  $+1$  thuộc về phía dương của siêu phẳng, tất cả các điểm được gán nhãn  $-1$  thuộc về phía âm của siêu phẳng).



Hình 2.2 Tập dữ liệu được phân chia tuyến tính

Để xác định hàm phân lớp dựa trên phương pháp SVM, ta sẽ tiến hành tìm hai siêu phẳng song song sao cho khoảng cách  $y$  giữa chúng là lớn nhất có thể để phân tách hai lớp này ra làm hai phía. Hàm phân tách tương ứng với phương trình siêu phẳng nằm giữa hai siêu phẳng tìm được.

### 2.1.3 Các bước chính của phương pháp SVM

- Tiền xử lý dữ liệu: Phương pháp SVM yêu cầu dữ liệu được diễn tả như các vector của các số thực. Như vậy nếu đầu vào chưa phải là số thực thì ta cần phải tìm cách chuyển chúng về dạng số của SVM. Tránh các số quá lớn, thường nên chuẩn hóa dữ liệu để chuyển về đoạn  $[-1,1]$  hoặc  $[0,1]$ .
- Chọn hàm nhân: Cần chọn hàm nhân phù hợp tương ứng cho từng bài toán toán cụ thể để đạt được độ chính xác cao trong quá trình phân lớp.
- Thực hiện việc kiểm tra chéo để xác định các tham số cho ứng dụng.
- Sử dụng các tham số cho việc huấn luyện tập mẫu.

- Kiểm thử tập dữ liệu Test.

### 2.1.4 Ưu điểm phương pháp SVM trong phân lớp dữ liệu

Chúng ta có thể thấy các thuật toán phân lớp hai lớp như SVM đều có đặc điểm chung là yêu cầu dữ liệu phải được biểu diễn dưới dạng vector đặc trưng, tuy nhiên các thuật toán khác đều phải sử dụng các ước lượng tham số và ngưỡng tối ưu trong khi đó thuật toán SVM có thể tự tìm ra các tham số tối ưu này. Trong các phương pháp thì SVM là phương pháp sử dụng không gian vector đặc trưng lớn nhất (hơn 10.000 chiều) trong khi đó các phương pháp khác có số chiều bé hơn nhiều (như Naïve Bayes là 2000, k-Nearest Neighbors là 2415...).

## 2.2 Giới thiệu về dữ liệu sử dụng

Dữ liệu được sử dụng trong luận văn đề dự đoán giới tính là tập dữ liệu có sẵn PAKDD'15 do Công ty Cổ phần FPT cung cấp ([www.fpt.com.vn](http://www.fpt.com.vn)). Dữ liệu được lấy từ lịch sử các hoạt động xem trang Web sản phẩm của người dùng. Nội dung dữ liệu đã thu thập được bao gồm các đặc trưng về thời gian và các danh mục chủng loại sản phẩm.

**Định dạng dữ liệu:** Dữ liệu đã cho được chia thành tập huấn luyện (*trainingData.csv*) và tập thử nghiệm (*testData.csv*) riêng biệt. Mỗi file này chứa 15.000 bản ghi tương ứng với nhật ký xem sản phẩm của người dùng Internet. Mỗi bản ghi hoạt động truy cập bao gồm bốn loại thông tin được phân cách bằng dấu phẩy (,). Loại thông tin đầu tiên là ID nhật ký hoạt động. Loại thông tin thứ hai và thứ ba tương ứng với thời gian bắt đầu truy cập và thời gian kết thúc truy cập. Thông tin cuối cùng là danh sách các danh mục, sản phẩm được người dùng truy cập trong một phiên hoạt động. Thứ tự truy cập danh mục, sản phẩm trong một phiên được phân cách nhau bởi dấu chấm phẩy (;). Ngoài ra còn có 2 file nhãn Label là *trainningLabel* và *testLabel* tương ứng với tập huấn luyện và tập thử nghiệm chứa hai loại thông tin về giới tính là nam và nữ.

Ví dụ minh họa trong tập dữ liệu, mỗi bản ghi gồm có các thông tin như sau:

- Session ID
- Start time (thời gian bắt đầu phiên)
- End time (thời gian kết thúc phiên)
- Danh sách các ID sản phẩm

Ví dụ về một bản ghi lịch sử truy cập:

u10008, 2014/11/17 19:20:06, 2014/11/17 19:21:54,

A00001 / B00001 / C00001 / D00001 /; A00001 / B00002 / C00002 / D00002

Từ tập dữ liệu cung cấp các thông tin ở trên, ta chia thông tin thành hai loại đặc trưng chính: Đặc trưng theo mốc thời gian và đặc trưng về các danh mục chủng loại sản phẩm.

- Đặc trưng về thời gian được đại diện bằng hai loại thông tin là thời gian bắt đầu truy cập và thời gian kết thúc truy cập bao gồm các thuộc tính như Ngày, tháng, năm, giờ, phút.
- Đặc trưng về danh mục, sản phẩm trong mỗi bản ghi nhật ký hoạt động được phân thành 4 cấp theo các chữ cái bắt đầu là A, B, C, D. Thứ tự các cấp được bảo toàn trong mỗi lượt xem. Mỗi bản ghi, danh mục lớn nhất bắt đầu bằng chữ cái A, các danh mục con tiếp theo bắt đầu bằng B, C và sản phẩm xem là D.

## **2.3 Các dạng đặc trưng dùng trong phân lớp**

### **2.3.1 Dạng đặc trưng theo mốc thời gian**

Trong tập dữ liệu PAKDD'15, đặc trưng về thời gian được biểu diễn thành hai loại thông tin là thời gian bắt đầu truy cập và thời gian kết thúc truy cập của người dùng mạng. Thời gian trong ngày, ngày trong tháng, tháng trong năm, thời gian bắt đầu xem, thời gian xem trong một bản ghi lịch sử truy cập, .. là những yếu tố có thể được sử dụng để dự đoán giới tính của một người dùng Internet.

Các đặc trưng theo mốc thời gian là quá trình theo dõi thời gian truy cập của người dùng mạng. Thông thường, thời gian truy cập của phiên hoạt động cao sẽ cho thấy người dùng đang mất nhiều thời gian hơn để xem các danh mục và chủng loại sản phẩm và các sản phẩm liên quan. Yếu tố này thường xảy ra với người dùng là nữ giới bởi các thói quen mua sắm và tìm hiểu thông tin của mình.

### **2.3.2 Dạng đặc trưng về danh mục và chủng loại sản phẩm**

Dữ liệu thu thập được của tập huấn luyện gồm có 15000 bản ghi nhật ký hoạt động truy cập. Các bản ghi này đều ghi lại quá trình truy cập sản phẩm của người dùng mạng. Quá trình này được bắt đầu từ danh mục cấp 1, tiếp đến là các danh mục con cấp 2, từ danh mục con cấp 2 chuyển đến danh mục con cấp 3 và từ danh mục con cấp 3 truy cập đến đích cuối là sản phẩm. Dựa vào tập dữ liệu và nhật ký hoạt động của người dùng, ta có thể phân cấp quá trình truy cập thành 4 loại đặc trưng của danh mục, sản phẩm theo các cấp tương ứng là A, B, C, D. Trong đó số lượng Danh mục A (cấp 1) là 11, số lượng danh mục B (cấp 2) là 86, số lượng danh mục C (cấp 3) là 383 và Sản phẩm D là 21881.

Bảng 2.1. Tóm tắt các đặc trưng dựa trên danh mục &amp; sản phẩm

Tên đặc trưng	Miêu tả
Danh mục A_ID	ID Danh mục cấp một
Danh mục B_ID	ID Danh mục cấp hai
Danh mục C_ID	ID Danh mục cấp ba
Sản phẩm D_ID	ID Sản phẩm theo các danh mục

Các danh mục chủng loại sản phẩm chứa các đặc trưng và yếu tố quan trọng cho việc phân loại giới tính của người dùng. Phần lớn các danh mục sản phẩm được người dùng truy cập thể hiện sự quan tâm, sở thích dựa trên các yếu tố về giới tính. Ví dụ như các sản phẩm liên quan đến thời trang, mỹ phẩm thường được xem bởi người dùng là giới tính nữ, các sản phẩm liên quan đến thể thao, công nghệ thì hay được truy cập bởi nam giới.

Từ tập dữ liệu đã cho, tác giả đã phân tích số lượt truy cập của nam giới và nữ giới đối với các danh mục, sản phẩm thu được từ các nhật ký xem sản phẩm của họ. Số liệu thống kê truy cập được phân chia theo từng cấp danh mục chủng loại sản phẩm được thể hiện trong hình 2.10.

Danh mục A ( cấp 1 )			
proid	female	male	total
A00001	557	1792	2349
A00002	8726	1109	9835
A00003	2210	228	2438
A00005	271	57	328
A00010	35	40	75
A00008	15	21	36
A00004	74	133	207
A00007	27	13	40
A00006	139	30	169
A00009	28	16	44
A00011	54	22	76

Danh mục B ( cấp 2 )			
proid	female	male	total
B00001	1544	775	2319
B00002	3553	336	3889
B00006	265	22	287
B00007	881	114	995
B00011	73	37	110
B00003	2030	267	2297
B00022	495	60	555
B00005	450	45	495
B00051	20	7	27
B00004	800	400	1200
B00044	31	21	52

Danh mục C ( cấp 3 )			
proid	female	male	total
C00001	23	132	155
C00002	854	84	938
C00012	186	117	303
C00011	47	7	54
C00004	335	56	391
C00003	813	74	887
C00007	1619	133	1752
C00015	152	13	165
C00016	349	40	389
C00030	132	12	144
C00180	17	6	23

Sản phẩm D ( cấp 4 )			
proid	female	male	total
D00001	1	1	2
D24897	5	0	5
D00009	1	0	1
D00007	1	0	1
D00003	1	0	1
D00002	1	0	1
D00017	1	0	1
D00010	1	0	1
D00014	1	0	1
D00011	1	0	1
D00012	1	0	1

Hình 2.1 Số liệu thống kê truy cập theo các cấp danh mục chủng loại sản phẩm

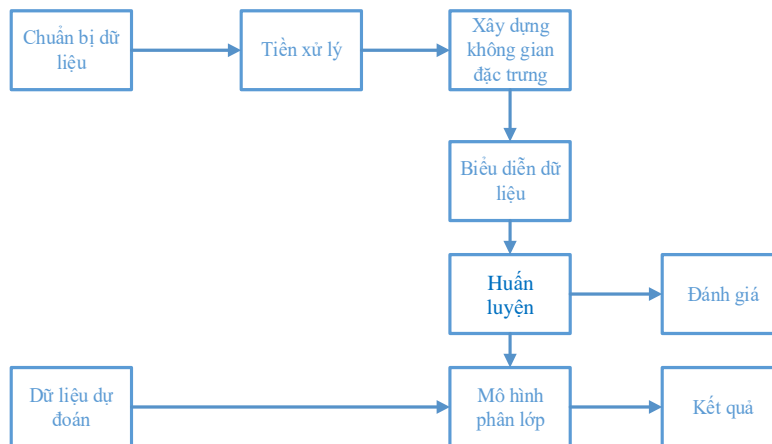
Từ các thống kê theo các cấp danh mục, sản phẩm trong tập dữ liệu cho thấy sự phân loại một cách rõ rệt giữa nam giới và nữ giới thông qua hoạt động truy cập các danh mục, sản phẩm dựa theo mức độ và tỉ lệ truy cập. Điều này cho thấy các đặc trưng về danh mục và chủng loại sản phẩm là loại đặc trưng chính trong việc xác định giới tính.

## 2.4 Xây dựng mô hình dự đoán giới tính dựa trên học máy có giám sát

Trong luận văn, với mục đích sử dụng phương pháp SVM để phân loại giới tính do đã được nhiều công trình đánh giá có độ chính xác cao trong phân lớp văn bản. Tuy nhiên, để đánh giá với đặc điểm của dạng dữ liệu ngắn, tác giả sẽ sử dụng phương pháp đếm số lần xuất



hiện của đặc trưng có trong từng bản ghi để phân loại, từ đó so sánh và rút ra kết luận cuối cùng về lựa chọn phương pháp phân loại cho bài toán.



**Hình 2.2 Mô hình phân loại dự đoán giới tính người dùng Internet**

#### 2.4.1 Tiền xử lý dữ liệu

Đây là giai đoạn "làm sạch" dữ liệu, bao gồm các bước sau:

- Loại bỏ trường không cần thiết
- Loại bỏ các ký tự đặc biệt ([ ], [, ], [:], [;], [/])
- Tách các đặc trưng

#### 2.4.2 Biểu diễn dữ liệu

Để có thể tiến hành thực nghiệm và đánh giá kết quả phân loại giới tính dựa trên tập dữ liệu lịch sử truy cập ta cần phải chuẩn hóa dữ liệu và áp dụng phương pháp cho tập dữ liệu đã có nhằm tiết kiệm không gian lưu trữ và gia tăng tốc độ xử lý.

Bước 1: Ta loại bỏ thuộc tính id vì thuộc tính này không dùng trong mô hình.

Bước 2: Tạo một danh sách các thuộc tính theo thứ tự bắt đầu từ các thuộc tính ngày, tháng, năm, giờ truy cập, phút sau đó là các thuộc tính phân cấp các danh mục và chủng loại sản phẩm A, B, C, D.

Bước 3: Định dạng thuộc tính thời gian và đếm số lần xuất hiện của các thuộc tính danh mục, sản phẩm có trong bản ghi (Giá trị, vị trí)

Thuộc tính thời gian	Thuộc tính danh mục, sp	Nhãn
7,1 12,2 2014,3 20,4 2,5 2,6 2,8 2,10 1,16 1,19		female
7,1 12,2 2014,3 20,4 0,5 2,6 2,8 1,10 1,11 1,15 1,17		female
20,1 12,2 2014,3 22,4 0,5 1,6 1,8 1,10 1,18		male
14,1 11,2 2014, 3 0,4 3,5 2,6 2,9 2,12 1,13 1,14		male

## 2.5 Kết luận chương

Chương này em đã giới thiệu chi tiết thuật toán Máy vector hỗ trợ SVM và giới thiệu về tập dữ liệu PAKDD'15 sử dụng trong luận văn. Ngoài ra, em mô tả cụ thể các dạng đặc trưng có trong dữ liệu. Dạng đặc trưng thứ nhất là dạng được trưng theo dấu thời gian truy cập của người dùng Internet. Dạng đặc trưng thứ hai là dạng đặc trưng theo danh mục và sản phẩm. Từ hai dạng đặc trưng này ta có thể khai phá và phân loại giới tính cho tập học và xây dựng mô hình dự đoán giới tính dựa trên phương pháp học máy SVM.

## CHƯƠNG 3: THỬ NGHIỆM VÀ ĐÁNH GIÁ

### 3.1 Mô tả dữ liệu

Để bắt đầu quá trình phân loại, ta cần xây dựng một tập huấn luyện theo đúng định dạng.

Sau các bước tiền xử lý, lịch sử truy cập được biểu diễn dưới dạng:

$\langle \text{label}_i \rangle \langle \text{index}_1 \rangle : \langle \text{value}_1 \rangle \langle \text{index}_2 \rangle : \langle \text{value}_2 \rangle \dots \langle \text{index}_n \rangle : \langle \text{value}_n \rangle$

Trong đó:

- $\text{label}_i$  là giá trị đích của tập huấn luyện. Đối với việc phân loại, nó là một số nguyên xác định một lớp, nhãn
- $\text{index}_i$  là một số nguyên bắt đầu từ 1. Cụ thể trong bài toán phân loại nó đại diện cho các đặc trưng.
- $\text{value}_i$  là một số thực. Giá trị này thể hiện mức độ liên quan của đặc trưng đối với một phân loại nằm trong khoảng  $[-1, 1]$ . Do các đặc trưng trong phân loại giới tính đều là đặc trưng nhị phân nên lúc huấn luyện giá trị này sẽ là 1.

### 3.2 Các tiêu chuẩn đánh giá

Các tiêu chí đánh giá hiệu quả dựa vào các kết quả thu được từ phân lớp của giải thuật:

- Số đúng dương (TP- True positive): số phần tử dương được phân loại dương
- Số sai âm (FN - False negative): số phần tử dương được phân loại âm
- Số đúng âm (TN- True negative): số phần tử âm được phân loại âm
- Số sai dương (FP - False positive): số phần tử âm được phân loại dương
- Độ chính xác dùng cho đo lường (Precision) =  $TP / (TP + FP)$
- Độ bao phủ (Recall) =  $TP / (TP + FN)$
- Độ đo F-Score =  $2 * \text{Precision} * \text{Recall} / (\text{Precision} + \text{Recall})$
- Độ chính xác dùng cho kết quả (Accuracy) =  $(TP + TN) / (TP + FP + TN + FN)$

### 3.3 Phương pháp thực nghiệm

Tập dữ liệu huấn luyện và kiểm thử dự đoán giới tính chứa các thông tin: Ngày, tháng, năm, giờ truy cập, thời gian xem, danh mục cấp A, danh mục con cấp B, danh mục con cấp C và sản phẩm D. Để định dạng dữ liệu, chúng ta cần biết LibSVM\_Tool và Weka học thế nào. Trong học máy nó thường được gọi là “Bộ thuộc tính”. Trong trường hợp phân loại giới

tính chúng ta xem mỗi danh mục sản phẩm và thời gian truy cập như một thuộc tính và được sắp xếp theo thứ tự bắt đầu từ 1 cho đến thuộc tính cuối cùng.

Tác giả khai thác những đặc điểm này và đếm số lần xuất hiện của thuộc tính trong mỗi bản ghi tương ứng với giới tính đã cho, sau đó đưa ra tập dữ liệu huấn luyện theo định dạng dữ liệu của hai bộ công cụ. Để đưa ra các tập dữ liệu đã được xử lý tác giả tạo 1 project Java tên là **Gender\_Prediction\_Network**.

Input: Là tập dữ liệu huấn luyện (*trainingData.csv*) và tập dữ liệu thử nghiệm (*TestData.csv*).

Output: Là file định dạng \*.arff (theo Weka) và file \*.txt (theo LibSVM\_Tool) có chứa tập dữ liệu huấn luyện và tập dữ liệu thử nghiệm kèm theo nhãn (xóa các dấu cách thừa, dấu phẩy “,”, dấu chấm phẩy “;” và dấu gạch ngang “/”) với mỗi dòng là một bản ghi lịch sử truy cập.

Sau khi xử lý dữ liệu ta bắt đầu tiến hành huấn luyện dữ liệu. Đầu tiên, để đánh giá hiệu quả của phương pháp SVM, tác giả sử dụng phương thức kiểm tra chéo (10-fold cross-validation) trên tập học cùng với công cụ **grid.py** trong LibSVM. Công cụ này sẽ tìm ra hai tham số tối ưu sao cho kết quả phân loại khi áp dụng hai tham số này sẽ đưa ra tỉ lệ cao nhất.

Bảng 3.1 Hai tham số tối ưu cho các mô hình huấn luyện

Mô hình\ Tham số	C	Gamma
A	2.0	0.03125
B	32.0	0.0078125
C	8.0	0.03125
D	8.0	0.03125
ALL	32.0	0.0078125

Hai tham số tối ưu này sẽ được kết hợp trong quá trình huấn luyện với Cross-Validation. Một tập con sẽ được giữ lại để làm tập dữ liệu kiểm tra, còn 9 tập còn lại sẽ được sử dụng để huấn luyện SVM, sau đó SVM này sẽ được dùng để dự đoán trên tập dữ liệu kiểm tra. Quá trình này sẽ được lặp đi lặp lại 10 lần sao cho tất cả các tập con đều sẽ được chọn làm tập dữ liệu kiểm tra. Trong quá trình thực nghiệm ta chia ra thành 4 mô hình rời rạc theo các cấp danh mục và chủng loại sản phẩm A, B, C, D và 1 mô hình chính lấy tên là ALL bao gồm tất cả các đặc trưng của tập dữ liệu để huấn luyện để so sánh kết quả phân loại giới tính giữa các mô hình đặc trưng rời rạc và mô hình đặc trưng tổng thể. Các tiêu chuẩn đánh giá sẽ

được tính trung bình từ các giá trị có được từ 10 lần lặp đó. Kết quả phân loại giới tính của tập dữ liệu lịch sử truy cập theo các mô hình được trình bày trong mục 3.4.

### 3.4 Kết quả thực nghiệm

Kết quả thực nghiệm sử dụng 5 mô hình và 4 tiêu chuẩn đánh giá để đưa ra hiệu quả mô hình học máy SVM cho việc phân loại giới tính. Kết quả thu được cho thấy khả năng phân loại có độ chính xác cao nhưng giảm dần với các mô hình rời rạc từ mô hình danh mục A cho đến mô hình chủng loại sản phẩm D. Lý do bởi vì độ nhiều dữ liệu khá lớn đối với các mô hình rời rạc, mô hình càng nhiều đặc tính thì độ nhiễu càng lớn. Kết quả cụ thể với các mô hình rời rạc được thu thập trong các Bảng 3.2, 3.3, 3.4 và 3.5.

Bảng 3.2 Kết quả thu được với mô hình A

NHÃN	SVM Với Mô hình A			
	Precision	Recall	F-Score	Accuracy
Nam	77.4 %	55.3 %	64.5 %	86.51 %
Nữ	88.2 %	95.4 %	91.7 %	
Weighted Avg	85.8 %	86.5 %	85.6 %	

Bảng 3.3 Kết quả thu được với mô hình B

NHÃN	SVM Với Mô hình B			
	Precision	Recall	F-Score	Accuracy
Nam	75.8 %	37.5 %	50.2 %	83.48 %
Nữ	84.4 %	96.6 %	90.1 %	
Weighted Avg	82.5 %	83.5 %	81.2 %	

Bảng 3.4 Kết quả thu được với mô hình C

NHÃN	SVM Với Mô hình C			
	Precision	Recall	F-Score	Accuracy
Nam	73.9 %	40.6 %	52.4 %	83.63 %
Nữ	85 %	95.9 %	90.1 %	
Weighted Avg	82.5 %	83.6 %	81.7 %	

Bảng 3.5 Kết quả thu được với mô hình D

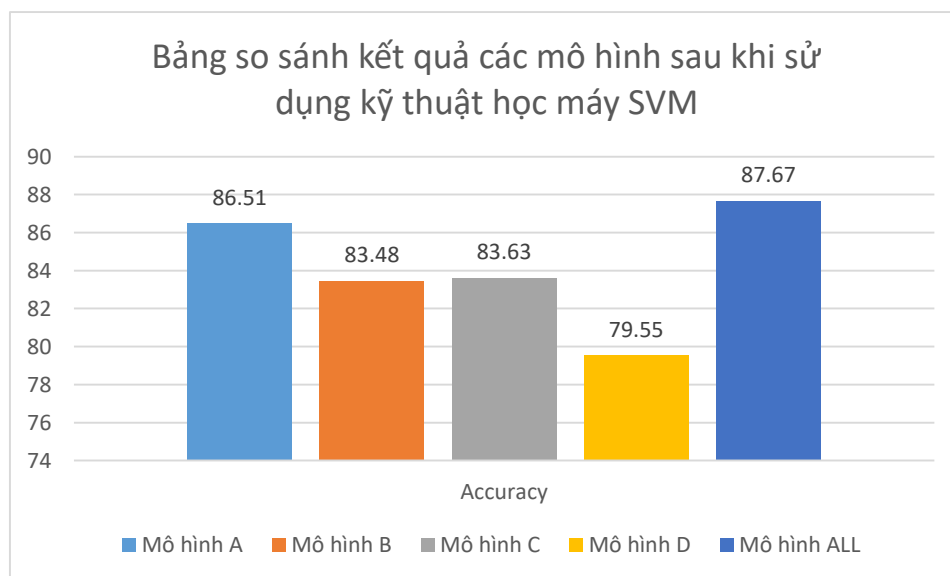
NHÃN	SVM Với Mô hình D			
	Precision	Recall	F-Score	Accuracy
Nam	82.7 %	9.9 %	17.7 %	79.55 %
Nữ	79.5 %	99.4 %	88.3 %	
Weighted Avg	80.2 %	79.5 %	72.6 %	

Tại Bảng 3.6 là bảng kết quả chính thu được từ mô hình học máy SVM khi sử dụng và kết hợp tất cả các đặc trưng và mô hình rời rạc của tập dữ liệu đã chuẩn hóa và đưa ra các tiêu chí đánh giá. So với kết quả của 4 mô hình rời rạc ở trên, tỉ lệ dự đoán khi kết hợp các đặc trưng lại với nhau mang đến tỉ lệ chính xác là 87.67 %. Từ các thực nghiệm trên cho thấy, SVM có độ phân lớp chính xác rất cao có thể đáp ứng được yêu cầu mà bài toán dự đoán giới tính đề ra.

Bảng 3.6 Kết quả thu được từ mô hình ALL

NHÃN	SVM với Mô hình All bao gồm tất cả các đặc trưng			
	Precision	Recall	F-Score	Accuracy
Nam	79.4 %	59.3 %	67.9 %	<b>87.67 %</b>
Nữ	89.3 %	95.7 %	92.4 %	
Weighted Avg	87.1 %	87.7 %	87 %	

Biểu đồ thể hiện độ chính xác của các mô hình:



### 3.5 So sánh với một số phương pháp khác

Để đánh giá thêm hiệu suất của mô hình dự đoán, luận văn đã tiến hành huấn luyện tập dữ liệu trên các mô hình học máy phổ biến khác là NaiveBayes và RandomTree, kết quả cụ thể được đưa ra trong bảng 3.7, 3.8.

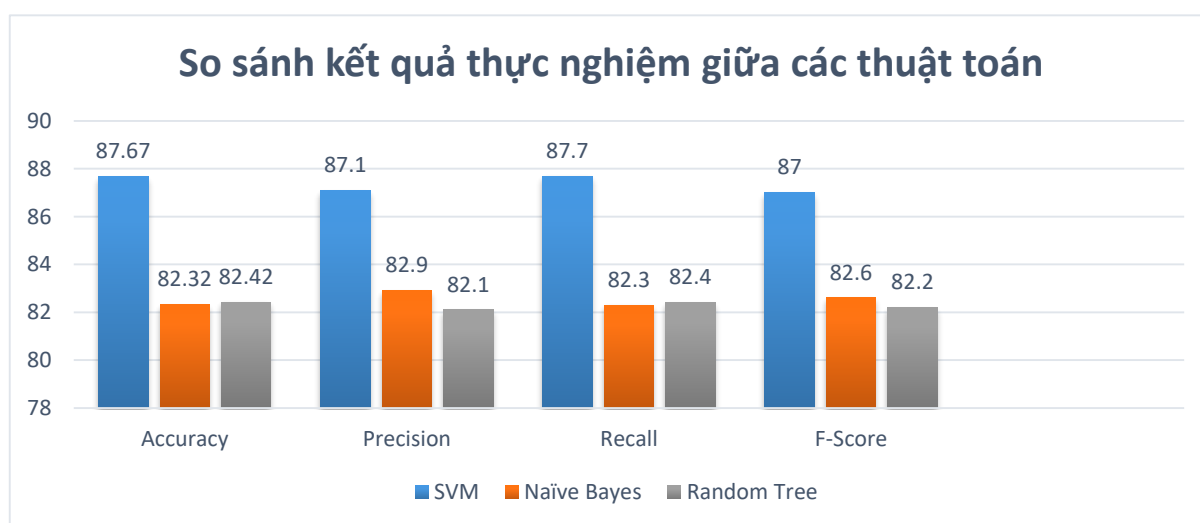
Bảng 3.7 Kết quả thu được từ mô hình Naïve Bayes

NHÃN	NaiveBayes			
	Precision	Recall	F-Score	Accuracy
Nam	59 %	64.3 %	61.5 %	82.32 %
Nữ	89.7 %	87.4 %	88.5 %	
Weighted Avg	82.9 %	82.3 %	82.6 %	

Bảng 3.8 Kết quả thu được từ mô hình Random Tree

NHÃN	Random Tree			
	Precision	Recall	F-Score	Accuracy
Nam	60.7 %	57 %	58.8 %	82.42 %
Nữ	88.1 %	89.6 %	88.8 %	
Weighted Avg	82.1 %	82.4 %	82.2 %	

Để dễ hình dung hơn thì chúng ta xem biểu đồ sau:



**Nhận xét:** Dựa vào bảng 3.6, 3.7, 3.8 tổng hợp kết quả phân loại giới tính trên các mô hình SVM, NaiveBayes, RandomTree ta nhận thấy được NaiveBayes cho kết quả thấp nhất khi phân loại mặc dù khả năng đưa ra độ chính xác khá cao với Accuracy = 82.32 % nhưng thực tế vẫn chưa tối ưu. Random Tree khá hơn nhưng tỉ lệ phân loại cũng chỉ nhiều hơn so với NaiveBayes là 0,1 %. Với SVM, tỉ lệ phân loại chính xác cao nhất so với 2 mô hình còn lại Accuracy = 87.67 %, ngoài ra các thông số Precision, Recall, F-Score cũng đều đưa ra tỉ lệ vượt trội. Kết quả này cho phép ta tin tưởng vào khả năng xử lý hiệu quả của mô hình học máy SVM cho vấn đề phân loại và xác định giới tính với dữ liệu có số chiều lớn.

### 3.6 Kết luận chương

Trong chương này, em đã nêu ra cách thức mô tả dữ liệu và chuẩn hóa dữ liệu Dữ liệu PAKDD'15 được sử dụng trong luận văn. Biểu diễn đặc trưng về danh mục sản phẩm và đặc trưng về thời gian truy cập của người dùng Internet để tạo ra dữ liệu huấn luyện để đưa vào các bộ công cụ hỗ trợ phân lớp cụ thể là LibSVM và Weka. Kết quả thực nghiệm được thể hiện trong 4 mô hình phân loại nhỏ và 1 mô hình phân loại tổng thể kết hợp với 4 tiêu chí đánh giá để đưa ra mức độ phù hợp của kỹ thuật học máy SVM khi áp dụng vào bài toán.

Do hạn chế về mặt thời gian, nên việc so sánh giữa các mô hình kỹ thuật học máy khác em chỉ đưa ra mô hình SVM với tất cả các đặc trưng và 2 mô hình huấn luyện là Naïve Bayes và Random Tree. Các kết quả thu được thể hiện trong Bảng 3.6, Bảng 3.7 và Bảng 3.8.

Kết quả thử nghiệm và đánh giá được tiến hành sau khi đã huấn luyện bộ dữ liệu theo 3 mô hình. Riêng trường hợp mô hình SVM thì có thêm công cụ grid.py giúp lựa chọn các tham số tối ưu. Kết quả thu được cho thấy SVM cho kết quả phân loại tốt hơn so với NaiveBayes và Random Tree với độ chính xác đạt trên 87 %.



## KẾT LUẬN

### 1. Kết quả đạt được

Luận văn tiến hành nghiên cứu giải quyết bài toán dự đoán giới tính người dùng Internet dựa trên lịch sử truy cập. Từ việc giải quyết bài toán giúp cho chúng ta tiến gần hơn đến sự thông minh của thế giới ảo, giúp quản lý tốt hơn hệ thống thông tin ngập tràn những nội dung. Bài toán là nền tảng cho nhiều ứng dụng quan trọng thực tế như quảng cáo nhắm mục tiêu, các hệ thống cung cấp tiếp thị dịch vụ thương mại tới đúng người dùng.

Những kết quả chính mà đồ án đạt được:

- Trình bày một cách khái quát, tổng quan nhất và nêu lên ý nghĩa, vai trò quan trọng của bài toán xác định giới tính người dùng Internet dựa trên lịch sử truy cập.
- Khảo sát nghiên cứu 3 phương pháp dự đoán giới tính đã có
- Đưa ra một số đặc trưng của tập dữ liệu lịch sử cho bài toán phân loại giới tính.
- Nghiên cứu và tìm hiểu về thuật toán Support Vector Machine trên hai lớp và nhiều lớp
- Nghiên cứu và làm thực nghiệm khi áp dụng Support Vector Machine để xác định giới tính của tập dữ liệu đã có.
- So sánh và phân tích các kết quả thực nghiệm với các mô hình thuật toán khác và đưa ra được trường hợp cho kết quả tốt nhất.

### 2. Hạn chế:

- Nghiên cứu dựa trên dữ liệu có sẵn, tập dữ liệu có sự mất cân bằng giới tính khi số lượng nữ nhiều hơn số lượng nam giới.
- Kết quả thực nghiệm đạt được vẫn chưa thực sự tốt so với kỳ vọng.
- Tốc độ xử lý dữ liệu vẫn chậm khi tập dữ liệu lớn

### 3. Hướng phát triển

- Thu thập bộ dữ liệu lớn hoàn chỉnh, phong phú về các lịch sử truy cập của người dùng Internet.
- Dựa trên nhiều đặc trưng để góp phần cải thiện khả năng phân loại và xác định giới tính người dùng áp dụng cho các bài toán thực tiễn
- Cải thiện hiệu suất, tăng tốc độ xử lý dữ liệu

- Ngoài ra em cũng sẽ nghiên cứu và thử nghiệm với một số mô hình thuật toán khác để tìm ra thuật toán phù hợp với bài toán xác định giới tính người dùng Internet.

### TÀI LIỆU THAM KHẢO

- [1] Do Viet Phuong and Tu Minh Phuong. “*Gender Prediction Using Browsing History*”. *KSE* (1) 2013: 271-283.
- [2] Hu, J., Zeng, H.-J., Li, H., Niu, C., Chen, Z. (2007) “*Demographic prediction based on user’s browsing behavior*”, Proceedings of the 16th international conference on World Wide Web, Banff, Alberta, Canada. [viewed 05.09.2016] Available from: <http://www.conference.org/www2007/papers/paper686.pdf>
- [3] Kabbur, S., Han, E.-H., Karypis, G. (2010) “*Content-based methods for predicting website demographic attributes*”, University of Minnesota Supercomputing Institute Research Report UMSI 2010/98 [viewed 06.09.2016] Available from: [http://www.dtc.umn.edu/publications/reports/2010\\_01.pdf](http://www.dtc.umn.edu/publications/reports/2010_01.pdf)
- [4] Speltdoorn, S. (2010) “*Predicting demographic characteristics of web users using semisupervised classification techniques*” Master’s dissertation, Ghent University, Faculty of Economics and Business Administration. [viewed 14.09.2016] Available from: [http://lib.ugent.be/fulltxt/RUG01/001/459/756/RUG01001459756\\_2011\\_0001\\_A\\_C.pdf](http://lib.ugent.be/fulltxt/RUG01/001/459/756/RUG01001459756_2011_0001_A_C.pdf)
- [5] Quanzeng You, Sumit Bhatia, Tong Sun, Jiebo Luo (2014) “*The eyes of the beholder: Gender prediction using images posted in Online Social Networks*”. Available from: [http://www.cs.rochester.edu/u/qyou/papers/gender\\_classification.pdf](http://www.cs.rochester.edu/u/qyou/papers/gender_classification.pdf)
- [6] Yuxiao Dong, Yang Yang, Jie Tang, Yang Yang, Nitesh V. Chawla (2014) “*Inferring User Demographics and Social Strategies in Mobile Social Networks*”. Available from: <http://www3.nd.edu/~ydong1/papers/KDD14-Dong-et-al-WhoAmI-demographic-prediction.pdf>
- [7] Yan, X., Yan, L.: Gender classification of weblogs authors. In: Proceedings of the AAAI Spring Symposium on Computational Approaches for Analyzing Weblogs, Stanford, CA, March 27-29, pp. 228–230 (2006). Available from: <http://aaai.org/Papers/Symposia/Spring/2006/SS-06-03/SS06-03-046.pdf>

- [8] Ying, J.J.C., Chang, Y.J., Huang, C.M., Tseng, V.S. (2012). Demographic prediction based on users mobile behaviors. Mobile Data Challenge. Available from: <http://www.idiap.ch/project/mdc/publications/files/mdc-final241ying.pdf>
- [9] Nguyen, D., Gravel, R., Trieschnigg, D., and Meder, T. (2013). "How old do you think i am?"; a study of language and age in twitter. Proceedings of the Seventh International AAAI Conference on Weblogs and Social Media. Available from: <http://www.dongnguyen.nl/publications/nguyen-icwsm2013.pdf>
- [10] Zhang, C., Zhang, P. (2010). Predicting gender from blog posts. Technical report, Technical Report. University of Massachusetts Amherst, USA.
- [11] Chang, C.C., Lin, C.J, 2001. LIBSVM – a library for support vector machines <http://www.csie.ntu.edu.tw/~cjlin/libsvm>
- [12] PENG Qiu-fang, LIU Yang – Research of gender prediciton based on SVM with E-commerce data. Available from: <http://lxbwk.njournal.sdu.edu.cn/EN/abstract/abstract3503.shtml>
- [13] Dong Nguyen, Rilana Gravel, Theo Meder, Dolf Trieschnigg – TweetGenie: Automatic Age Prediction From Tweets. Available from: <http://dolf.trieschnigg.nl/papers/SIGWEB.2013.nguyen.pdf>
- [14] Josh Jia-Ching **Ying**, Yao-Jen Chang, Chi-Min Huang and Vincent S. Tseng (2012) – Demographic Prediction Based on User's Mobile Behaviors. Available from: <http://www.idiap.ch/project/mdc/publications/files/mdc-final241-ying.pdf>
- [15] Zhang, C., Zhang, P. (2010) – Predicting gender from blog posts. Technical report, Technical Report. University of Massachusetts Amherst, USA.
- [16] Mendenhall's studies of word-length distribution in the works of Shakespeare and Bacon. Available from: <https://academic.oup.com/biomet/article-abstract/62/1/207/220350/Mendenhall-s-studies-of-word-length-distribution>
- [17] De Vel, O., Anderson, A., Corney, M., Mohay, G. M. (2001). Mining e-mail content for author identification forensics. SIGMOD Record 30(4), pp. 55-64.

- [18] Argamon, S., Koppel, M., Fine, J. and Shimoni, A. (2003). Gender, Genre, and Writing Style in Formal Written Texts, *Text* 23(3), August.
- [19] Argamon, S., Koppel, M., Pennebaker, J. and Schler, J. (2008). Automatically Profiling the Author of an Anonymous Text, *Communications of the ACM*.
- [20] Making Large-Scale SVM Learning Practical - Thorsten Joachims. Available from: [https://www.cs.cornell.edu/people/tj/publications/joachims\\_99a.pdf](https://www.cs.cornell.edu/people/tj/publications/joachims_99a.pdf)
- [21] Weka - Available from: <http://www.cs.waikato.ac.nz/ml/weka/>