

Lời cam đoan

Lời cảm ơn

Mục lục

Danh sách hình vẽ	5
Danh sách bảng	6
Danh sách từ viết tắt	7
1 Tổng quan	8
1.1 Giới thiệu đề tài	8
1.2 Mục tiêu và phạm vi đề tài	8
1.3 Cấu trúc luận văn	8
2 Các công trình liên quan	9
2.1 Bệnh án điện tử	9
2.2 Các hướng khai thác dữ liệu bệnh án điện tử	9
3 Kiến thức nền tảng	10
3.1 Các định nghĩa và thuật ngữ	10
3.2 Học máy	10
3.3 Support Vector Machine	10
3.4 Xử lý ngôn ngữ tự nhiên	10
3.5 Phân giải đồng tham chiếu	10
4 Phân tích yêu cầu và thiết kế tổng quát	11
4.1 Nội dung bài toán	11
4.2 Kiến trúc hệ thống	11
5 Chi tiết hệ thống	14
5.1 Tiền xử lý	14
5.2 Xây dựng cặp khái niệm	14
5.3 Rút trích đặc trưng	14
5.4 Gom cụm và xây dựng chuỗi đồng tham chiếu	17
5.5 Đánh giá hiệu năng	17
6 Thí nghiệm đánh giá	18
6.1 Tập dữ liệu	18
6.2 Kết quả	18

MỤC LỤC

7	Tổng kết	19
	Tài liệu tham khảo	19

Danh sách hình vẽ

4.1	Sơ đồ huấn luyện	12
4.2	Sơ đồ phân giải đồng tham chiếu	13
5.1	Tổng quan hệ thống phân giải đồng tham chiếu	15

Danh sách bảng

5.1	Tập đặc trưng cho lớp Person	16
-----	--	----

Danh sách từ viết tắt

BADT Bệnh án điện tử

Chương 1

Tổng quan

1.1 Giới thiệu đề tài

1.2 Mục tiêu và phạm vi đề tài

1.3 Cấu trúc luận văn

Chương 2

Các công trình liên quan

2.1 Bệnh án điện tử

2.2 Các hướng khai thác dữ liệu bệnh án điện tử

Chương 3

Kiến thức nền tảng

3.1 Các định nghĩa và thuật ngữ

3.2 Học máy

3.3 Support Vector Machine

3.4 Xử lý ngôn ngữ tự nhiên

3.5 Phân giải đồng tham chiếu

Chương 4

Phân tích yêu cầu và thiết kế tổng quát

4.1 Nội dung bài toán

4.2 Kiến trúc hệ thống

Dựa vào hệ thống có hiệu năng tốt nhất của thử thách i2b2 năm 2011 (hệ thống I), mô hình phân giải đồng tham chiếu mà chúng tôi sử dụng để hiện thực hệ thống là mô hình cặp thực thể. Tư tưởng cơ bản của mô hình này là xác định xem hai khái niệm bất kì có đồng tham chiếu với nhau hay không, sau đó gom nhóm các cặp đồng tham chiếu có một khái niệm chung lại để tạo thành các chuỗi đồng tham chiếu. Như vậy kiến trúc tổng quát của hệ thống chúng tôi hiện thực gồm 2 quy trình: *quy trình huấn luyện hệ thống phân loại* và *quy trình phân giải đồng tham chiếu*. Trong đó quy trình huấn luyện là bước huấn luyện các model phân loại dựa trên dữ liệu mẫu đã được phân giải đồng tham chiếu. Quy trình phân giải sử dụng các model phân loại đã được huấn luyện để xác định tính đồng tham chiếu của các cặp khái niệm, từ đó sử dụng một giải thuật gom nhóm các cặp đồng tham chiếu lại để tạo thành các chuỗi đồng tham chiếu.

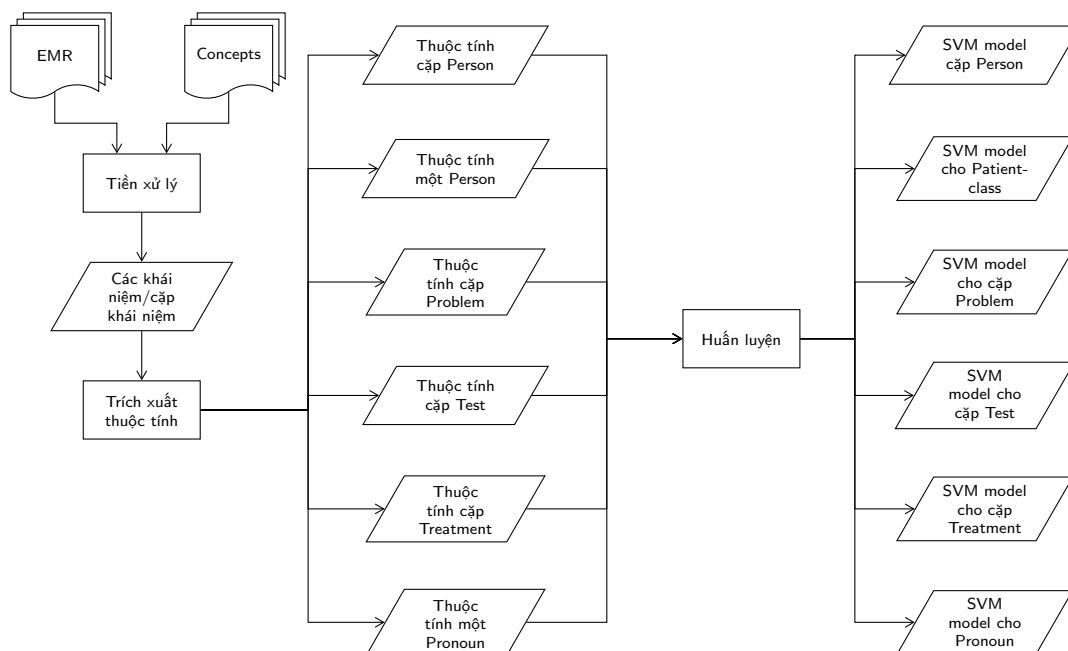
Quy trình huấn luyện hệ thống phân loại

Để xác định tính đồng tham chiếu giữa hai khái niệm bất kì, ta cần huấn luyện một model phân loại dựa trên dữ liệu mẫu. Vì đầu vào của quy trình là các văn bản BADT và danh sách các khái niệm đã được gán nhãn, hệ thống cần trích xuất thuộc tính các dữ liệu thô này rồi mới có thể đưa vào để huấn luyện. Bên cạnh đó, các khái niệm đã được phân loại vào 4 nhóm chính là Person, Problem, Test và Treatment, còn các đại từ được phân vào nhóm Pronoun nên để giảm bớt số cặp khái niệm được sinh ra, chúng tôi huấn luyện 4 model để xác định tính đồng tham chiếu của riêng các cặp Person-Person, Problem-Problem, Test-Test và Treatment-Treatment (vì hai khái niệm thuộc hai lớp khác nhau thì nghiêm nhiên không đồng tham chiếu với nhau). Đối với các đại từ thì thường chỉ tới một khái niệm ở trước đó, nên việc xác định xem một đại từ thực chất mang ý nghĩa của lớp nào trong 4 lớp chính Person, Problem, Test, Treatment là một việc quan trọng. Sau khi đã xác định được lớp chính, chúng tôi sẽ xem đại từ đang xét là

đồng tham chiếu với khái niệm thuộc lớp tương ứng ở gần nhất trước đó. Các ý này đều là của các tác giả hệ thống I.

Ngoài ra cũng theo các tác giả này, thông tin một khái niệm lớp Person có chỉ về bệnh nhân hay không góp một phần quan trọng trong việc phân loại đúng tính đồng tham chiếu của cặp các khái niệm lớp này. Trong miền văn bản BADT, các khái niệm chỉ người thường chỉ đề cập đến một trong ba loại: bệnh nhân, người thân của bệnh nhân và nhân sự của bệnh viện. Do một BADT, mà cụ thể là hồ sơ xuất viện, thông thường chỉ đề cập đến một bệnh nhân nên những khái niệm nào chỉ về bệnh nhân thì thường chắc chắn nằm trong cùng một chuỗi đồng tham chiếu lớn nhất và duy nhất chỉ về bệnh nhân đó. Từ nhận định này, nhóm tác giả của hệ thống I đã đưa vào thuộc tính Patient-class cho cặp hai khái niệm lớp Person, nó mang giá trị 1 khi hai khái niệm đều chỉ về bệnh nhân và 0 trong các trường hợp khác. Ở bước huấn luyện, thông tin “một khái niệm Person có chỉ về bệnh nhân hay không” được lấy từ tập chuỗi kết quả (ground truth), còn ở bước phân giải đồng tham chiếu thông tin này được xác định nhờ một model phân loại đã được huấn luyện.

Như vậy mục đích của quy trình huấn luyện là xây dựng tổng cộng 6 SVM model, trong đó 4 SVM model nhằm mục đích phân loại và đánh giá độ tin cậy đồng tham chiếu của các cặp khái niệm Person-Person, Problem-Problem, Test-Test và Treatment-Treatment; 1 SVM model để xác định các khái niệm Person có là bệnh nhân hay không (Patient-class) và 1 SVM model để phân loại các đại từ (các khái niệm lớp Pronoun) vào một trong bốn lớp Person, Problem, Test và Treatment. Đầu vào của quy trình này là toàn bộ các văn bản BADT với các khái niệm đã được trích xuất và gán nhãn. Sau khi tiền xử lý, hệ thống xây dựng các mẫu huấn luyện bao gồm: Person, Person-Person, Problem-Problem, Test-Test, Treatment-Treatment và Pronoun từ danh sách các khái niệm. Sáu tập mẫu này được trích xuất thuộc tính và đưa vào để huấn luyện 6 SVM model (Hình 4.1). Thư viện SVM được nhóm sử dụng là LibSVM.

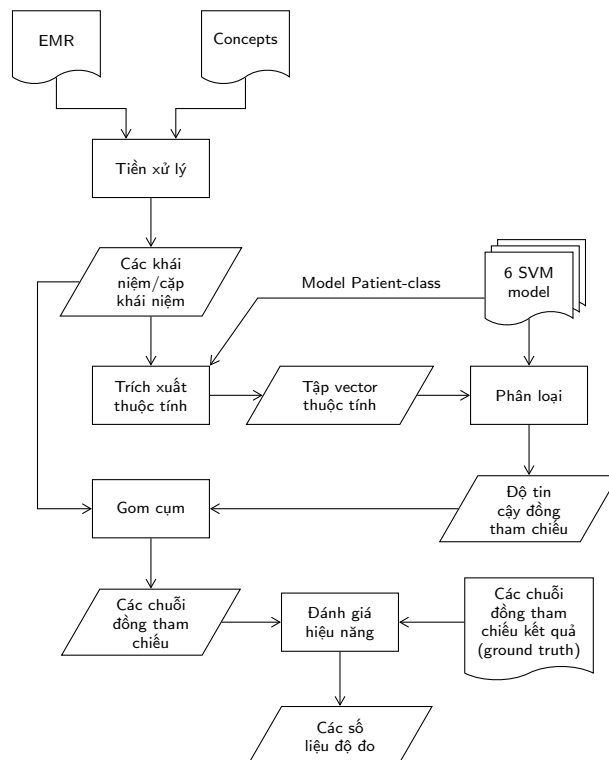


Hình 4.1: Sơ đồ huấn luyện

Quy trình phân giải đồng tham chiếu

Quy trình phân giải đồng tham chiếu sử dụng 6 SVM model đã được huấn luyện ở trên, cùng với đó là một giải thuật gom nhóm các cặp khái niệm đã được phân loại là đồng tham chiếu với nhau lại để cuối cùng tạo thành các chuỗi đồng tham chiếu. Có thể xem đây là quy trình mang đi ứng dụng thực tế để phân giải cho những văn bản BADT mới. Dựa vào hệ thống I, chúng tôi sử dụng giải thuật gom cụm tốt nhất trước để lựa chọn các cặp đồng tham chiếu có độ tin cậy cao nhất, sau đó xây dựng các chuỗi đồng tham chiếu bằng cách nối các cặp có một khái niệm chung. Đối với lớp Pronoun, sau khi đã xác định được lớp chính của một đại từ bất kì, chúng tôi tạo một cặp đồng tham chiếu giữa đại từ đó và khái niệm thuộc lớp chính tương ứng ở gần nhất trước đó trong văn bản. Theo nhận định của các tác giả hệ thống I, tuy cách làm này đơn giản nhưng lại tỏ ra rất hiệu quả.

Hình 4.2 mô tả trực quan quy trình phân giải đồng tham chiếu. Ở bước trích xuất thuộc tính của các cặp Person, chúng tôi sử dụng model phân loại bệnh nhân để xác định giá trị cho thuộc tính Patient-class đã được đề cập ở trên. Theo kết quả đánh giá các hệ thống dự thi thử thách i2b2 2011, ba hệ đo được sử dụng là: MUC, B-CUBED và CEAF. Chúng tôi cũng hiện thực các hệ đo này để đánh giá hệ thống của mình bằng cách so sánh với kết quả của hệ thống I. Vì các hệ đo này đánh giá trên các chuỗi đồng tham chiếu chứ không phải đánh giá hiệu năng phân loại thông thường nên cách thức hiện thực chúng có phần phức tạp, chúng tôi sẽ trình bày chi tiết các cách hiện thực này ở phần sau.



Hình 4.2: Sơ đồ phân giải đồng tham chiếu

Chương 5

Chi tiết hệ thống

5.1 Tiền xử lý

Trong quá trình rút trích đặc trưng, một số khái niệm được miêu tả cụ thể làm cho việc so trùng chuỗi hoặc tìm kiếm từ các nguồn tri thức nhân loại thiếu chính xác [1]. Ví dụ như khái niệm “her CT scan” và khái niệm “a CT scan”. Mặc dù hai khái niệm này cùng chỉ một thủ tục y tế nhưng không trùng chuỗi. Ngoài ra các mạo từ “her”, “a” làm việc tìm kiếm tri thức nhân loại từ các nguồn tri thức như Wikipedia, WordNet không được chính xác hoặc không thể tìm được kết quả. Vì vậy trước khi rút trích đặc trưng, các khái niệm được tiền xử lý để loại bỏ mạo từ và các thông tin ngữ cảnh. Tuy nhiên, quá trình tiền xử lý chỉ được áp dụng cho các đặc trưng liên quan so trùng chuỗi và tìm kiếm tri thức nhân loại, các đặc trưng khác không cần qua quá trình tiền xử lý mà nhận vào nguyên gốc khái niệm được xác định.

Quá trình tiền xử lý gồm hai bước. Đầu tiên khái niệm sẽ được loại bỏ tất cả mạo từ. Sau đó, nếu khái niệm có bao gồm giới từ thì giới từ đó và toàn bộ nội dung theo sau sẽ được lược bỏ. Ví dụ như khái niệm “an MRI of the knee” sau quá trình tiền xử lý sẽ trở thành “MRI”. Danh sách mạo từ được xây dựng từ tập dữ liệu và các mạo từ thông dụng của tiếng Anh.

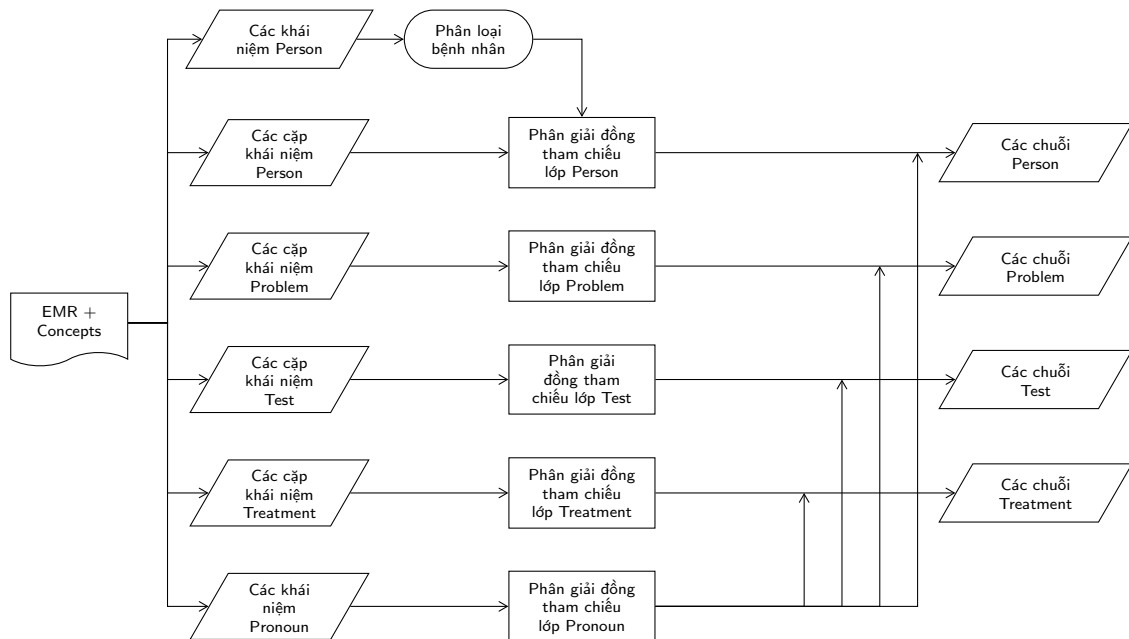
Đặc biệt các khái niệm thuộc lớp Problem/Treatment/Test thường được kèm thêm thông tin về định lượng như 10mg, 5 lit và các thông tin về vị trí giải phẫu học như “upper”, “left”, “right”. Để tăng khả năng tìm kiếm tri thức nhân loại, chúng tôi đề xuất loại bỏ các thông tin ngữ cảnh về số, định lượng và vị trí giải phẫu khỏi khái niệm. Các thông tin ngữ cảnh được loại bỏ bằng cách sử dụng biểu thức chính quy và các từ vựng được xây dựng từ tập dữ liệu. Các đặc trưng liên quan so trùng chuỗi không áp dụng bước tiền xử lý loại bỏ thông tin ngữ cảnh này.

5.2 Xây dựng cặp khái niệm

5.3 Rút trích đặc trưng

Từ các phân tích được đề cập ở Phần 3, ngoài các thuộc tính chung về mặt ngôn ngữ (như ngữ pháp hay từ vựng), từng lớp khái niệm ở BADT còn mang những đặc tính khác nhau. Việc này đòi hỏi chúng tôi phải thiết kế ba hệ thống rút trích đặc trưng và phân loại tương ứng khác nhau cho lớp Person, lớp Problem/Treatment/Test và lớp Pronoun. Hình 5.1 mô tả tổng quan

các hệ thống này, trong đó các khối “Đồng tham chiếu lớp X” bao hàm cả Hệ thống rút trích đặc trưng và Hệ thống phân loại cho lớp tương ứng.



Hình 5.1: Tổng quan hệ thống phân giải đồng tham chiếu

Nhóm Person

Tổng quát, các khái niệm thuộc lớp Person có thể là các đại từ nhân xưng (he, she, it, they, ...), đại từ sở hữu (his, her, its, their, ...), đại từ phản thân (himself, herself, itself, themselves, ...) hoặc tên người (Stephanie I Sept, Mr. Anders, Heidi Laura Md, ...). Việc phân giải đồng tham chiếu cho tên người và đại từ là công việc khó, vì thông tin có được từ các đại từ và tên người là rất ít. Ngoài ra trong một văn bản thường đề cập đến nhiều hơn một người, khiến cho việc phát hiện chính xác chuỗi đồng tham chiếu cho các khái niệm này là một thách thức lớn.

Việc giới hạn vấn đề lại trong phạm vi BADT giúp công việc này trở nên đơn giản hơn. Trong BADT, các khái niệm thuộc lớp Person thường được chia vào ba nhóm chính: bệnh nhân, người thân của bệnh nhân hoặc nhân sự của bệnh viện. Trong đó bệnh nhân là nhóm có số lượng khái niệm được đề cập nhiều nhất và chiếm phần lớn tổng số khái niệm lớp Person. Do vậy việc xác định một khái niệm thuộc vào nhóm nào đóng vai trò quan trọng trong việc phân giải chính xác chuỗi đồng tham chiếu cho khái niệm đó [1]. Từ lí do trên, đặc trưng có phải là bệnh nhân hay không được thêm vào hệ thống. Đặc trưng lớp Patient được xác định bằng phương pháp phân loại nhị phân SVM. Hai nhóm người thân của bệnh nhân và nhân sự của bệnh viện được xác định bằng các đặc trưng từ vựng. Bảng 5.1 trình bày đầy đủ các đặc trưng dùng cho lớp Person

Với các đặc trưng Name match, Relative match, Department match, Doctor title match, Doctor general match, Twin/Triplet, We, You, I, Pronoun match, chúng tôi hiện thực bằng cách xây dựng tập từ điển tương ứng với từng đặc trưng dựa trên việc khảo sát tập dữ liệu và sử dụng các biểu thức chính quy.

Bảng 5.1: Tập đặc trưng cho lớp Person

Thuộc tính	Giá trị	Giải thích
Patient-class	0, 1, 2	Không có khái niệm nào là bệnh nhân (0), cả hai khái niệm đều là bệnh nhân (1), trường hợp khác (2)
Distance between sentences	0, 1, 2, 3, ...	Số câu xuất hiện giữa hai khái niệm được xét
Distance between mentions	0, 1, 2, 3, ...	Số khái niệm xuất hiện giữa hai khái niệm được xét
String match	0, 1	Trùng chuỗi hoàn toàn (1), ngược lại (0)
Levenshtein distance between two mentions	(0, 1)	Khoảng cách Levenshtein giữa hai khái niệm
Number	0, 1, 2	Cả hai đều là số ít hoặc nhiều (1), ngược lại (0), không xác định (2)
Gender	0, 1, 2	Cùng giới tính (1), khác giới tính (0), không xác định (2)
Apposition	0, 1	Là đồng vị ngữ (1), ngược lại (0)
Alias	0, 1	Là từ viết tắt hoặc cùng nghĩa (1), ngược lại (0)
Who	0, 1	Nếu hai khái niệm liên kế nhau và được phân cách bởi dấu “:”
Name match	0, 1	Loại bỏ các “stop word” (dr, dr., mr, ...), so trùng chuỗi con, trùng (1), không trùng (0)
Relative match	0, 1	Cả hai đều cùng chỉ đến một thân nhân (1), ngược lại (0)
Department match	0, 1	Cả hai cùng chỉ đến một lĩnh vực y học (1), ngược lại (0)
Doctor title match	0, 1	Cả hai có cùng một chức vụ bác sĩ (1), ngược lại (0)
Doctor general match	0, 1	Cả hai cùng đề cập đến bác sĩ nói chung (1), ngược lại (0)
Twin/triplet	0, 1	Cả hai đều chỉ về cùng cặp sinh đôi/sinh ba (1), ngược lại (0)
We	0, 1	Cả hai đều chứa thông tin về “chúng tôi” (1), ngược lại (0)
You	0, 1	Cả hai đều chứa thông tin về “tôi” (1), ngược lại (0)
I	0, 1	Cả hai đều chứa thông tin về “bạn” (1), ngược lại (0)
Pronoun match	0, 1	Cả hai đều là đại từ chỉ người (1), ngược lại (0)

Đặc trưng về Giới tính được chúng tôi xác định dựa trên ba bước phân loại [2]. Bước thứ nhất: kiểm tra khái niệm có chứa các đại từ xác định giới tính như “Mr”, “Ms”, “she”, “he”, ... hay không. Nếu có, xác định giới tính dựa trên đại từ xuất hiện. Nếu không thực hiện bước thứ hai: kiểm tra khái niệm có xuất hiện nhiều hơn một lần hay không. Nếu xuất hiện nhiều hơn một lần thì các lần xuất hiện có chứa đại từ xác định giới tính hay không. Ví dụ khái niệm “Peter H. Diller” có thể xuất hiện nhiều lần, trong đó có xuất hiện dưới hình thức “Mr. Diller”. Nếu không thể xác định giới tính qua hai bước kiểm tra, khái niệm sẽ được phân loại bằng cách sử dụng cơ sở dữ liệu về tên tiếng Anh của hệ thống Apache OpenNLP.

Nhóm Patient

Nhóm Pronoun

Nhóm Problem/Test/Treatment

5.4 Gom cụm và xây dựng chuỗi đồng tham chiếu

5.5 Đánh giá hiệu năng

Hệ đo MUC

Hệ đo B-CUBED

Hệ đo CEAF

Chương 6

Thí nghiệm đánh giá

6.1 Tập dữ liệu

6.2 Kết quả

Chương 7

Tổng kết

Tài liệu tham khảo

- [1] Y. Xu, “A classification approach to coreference in discharge summaries: 2011 i2b2 challenge,” *Journal of the American Medical Informatics Association : JAMIA*, vol. 19, no. 5, pp. 897–905, 2012.
- [2] W. M. Soon, H. T. Ng, and D. C. Y. Lim, “A machine learning approach to coreference resolution of noun phrases,” *Computational Linguistics*, vol. 27, pp. 521–544, December 2001.