

ĐẠI HỌC QUỐC GIA TP.HCM
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN



KHOA KHOA HỌC MÁY TÍNH

CS221.O22 - XỬ LÝ NGÔN NGỮ TỰ NHIÊN

BÁO CÁO ĐỒ ÁN CUỐI KÌ

MULTI-STAGE DOCUMENT RANKING
FOR VIETNAMESE NEWS RETRIEVAL

GV hướng dẫn: TS. Nguyễn Trọng Chính, ThS. Nguyễn Đức Vũ,
ThS. Đặng Văn Thìn

Nhóm thực hiện

21521992 - Võ Đức Dương

21521465 - Trần Ngọc Thiện

21520617 - Nguyễn Thúc Bảo



Mục lục

1	Lý do lựa chọn đề tài	2
2	Giới thiệu bài toán	3
2.1	Mô tả bài toán	3
2.2	Mô tả Input/Output	3
2.3	Mô tả ngữ liệu	3
3	Giới thiệu phương pháp chính	5
3.1	Quy trình xử lý của Multi-stage Document Ranking	5
3.2	Các phương pháp term-based	6
3.2.1	TF-IDF	6
3.2.2	BM25	7
3.3	Pre-trained language model	8
3.3.1	PhoBERT	8
3.3.2	Vietnamese-SBERT	8
4	Cài đặt phương pháp	9
4.1	Tiền xử lý dữ liệu	9
4.2	Đánh giá bài báo bằng Cosine Similarity	9
5	Đánh giá phương pháp	10
5.1	Tập đánh giá	11
5.2	Kết quả thu được	12
6	Tài liệu tham khảo	13



1 Lý do lựa chọn đề tài

Trong thời đại cách mạng công nghiệp 4.0, Internet trở thành nguồn tài nguyên thông tin khổng lồ. Việc sử dụng các công cụ tìm kiếm như Google, Bing, ... để tra cứu thông tin trở nên phổ biến và quen thuộc với con người. Thế nhưng, **làm thế nào mà các công cụ tìm kiếm này có thể hiểu được câu truy vấn và trả về các thông tin liên quan?** Bắt nguồn từ câu hỏi này, nhóm quyết định tìm hiểu để xây dựng và tái hiện lại một hệ thống tương tự.

Theo ý tưởng đó, nhóm tiếp tục xem xét: **Đâu là thông tin mà người đọc thường sẽ chú ý nhiều hơn?** Trên thực tế, với mỗi một câu truy vấn chứa thông tin cần tra cứu, công cụ tìm kiếm có thể trả về vô số bài báo, bài viết với vô số đường dẫn chứa thông tin liên quan. Mỗi một câu truy vấn sẽ yêu cầu thông tin về một lĩnh vực, và ứng với mỗi lĩnh vực đó, những bài viết hoặc trang web ưu tiên trả về sẽ khác nhau.

Khi tìm kiếm *Thông tin hôm nay*, 10 kết quả trả về đầu tiên đều là các trang web báo điện tử. Lý giải điều này, nhóm cho rằng đây là các trang báo điện tử phổ biến ở Việt Nam, cập nhật các thông tin mới nhất theo từng ngày, dẫn đến số lượng truy cập đông đảo. Cùng với các dòng văn bản như *mới nhất 24h hôm nay, tin nhanh, tin nóng, Thời sự hôm nay, ...* đã định hướng để trình tìm kiếm trả về kết quả như trên.

Trong phạm vi đồ án môn học, nhóm chỉ thực hiện trong giới hạn là các thông tin văn bản, cụ thể hơn là các bài báo điện tử vì các trang báo điện tử luôn được người đọc đón nhận và truy cập hằng ngày. Tất cả những điều trên là lý do để nhóm thực hiện đề tài.

Đồ án này thể hiện những gì mà chúng em học được từ lớp **CS221.O22**. Chúng em xin gửi lời cảm ơn **TS. Nguyễn Trọng Chính, ThS. Nguyễn Đức Vũ, ThS. Đặng Văn Thìn** đã nhiệt tình giảng dạy, giải đáp thắc mắc trong quá trình học tập và thực hiện đồ án cuối kỳ. Đồ án chắc chắn không tránh khỏi những sai sót, chúng em hi vọng nhận được những góp ý của các thầy về đồ án để cải thiện và hoàn chỉnh hơn.



2 Giới thiệu bài toán

2.1 Mô tả bài toán

Cho một truy vấn q và một ngữ liệu D là tập hợp nhiều bài báo khác nhau. Bài toán đặt ra là xây dựng một hệ thống nhận vào truy vấn q và trả về các bài báo từ D , trong đó:

- Chứa các ngữ cảnh quan trọng liên quan đến q .
- Các bài báo được xếp hạng từ trên xuống dưới theo **mức độ liên quan** đến truy vấn. Các bài báo có độ liên quan cao hơn thì có xếp hạng (rank) cao hơn.

2.2 Mô tả Input/Output

a) Input:

- 1 tập ngữ liệu D : chứa thông tin các bài báo trên các trang báo điện tử. Thông tin chi tiết về tập ngữ liệu sẽ được mô tả ở mục **2.3**
- 1 câu truy vấn q bằng Tiếng Việt, chứa thông tin về bài báo cần tìm kiếm

b) Output: Một danh sách các bài báo trong ngữ liệu D , được sắp xếp giảm dần về mức độ liên quan của bài báo với truy vấn.

2.3 Mô tả ngữ liệu

Tập ngữ liệu chứa thông tin về các bài báo được thu thập từ 5 trang báo điện tử của Việt Nam, bao gồm: báo Lao Động, báo Dân Trí, báo VnExpresses, báo VTC và báo Đảng Cộng Sản. Với mỗi bài báo, ta thu thập 5 đặc trưng (feature) bao gồm: **title** (tiêu đề), **abstract** (tổng quan), **source** (nguồn bài báo), **link** (đường dẫn bài báo) và **topic** (chủ đề).



Điểm thi Toán, Văn, Ngoại ngữ vào lớp 10 của hơn 105.000 thí sinh được công bố lúc 17h.

Hình 1: Một bài báo Giáo dục từ báo VnExpress

Ví dụ, với một bài báo như *Hình 1*, ta có các thông tin sau:

- **title:** Hà Nội công bố điểm thi lớp 10
- **abstract:** Điểm thi Toán, Văn, Ngoại ngữ vào lớp 10 của hơn 105.000 thí sinh được công bố lúc 17h.
- **source:** Báo VnExpress
- **link:** <https://vnexpress.net/tra-cuu-diem-thi-lop-10-ha-noi-nam-2024-4762421.html>
- **topic:** Giáo dục

Tổng số lượng thu thập là 48357 bài báo. Ở đây, nhóm chỉ tập trung vào mức độ liên quan của tiêu đề và tổng quan đối với truy vấn nên thời gian của bài báo sẽ không được thu thập.

title	abstract	source	link	topic
Chợ gạo đã "khô" ở thuế thu nhập cá nhân	(ĐCSVN) – Có nhiều ý kiến cho rằng, mức giảm thuế giá trị gia tăng trong cách	Báo Đảng Cộng Sản	https://dangcongson.vn/cung-ban-luan/chua-giau-da-kho-ve-thue-thu-nhap-ca-nhan	Công bản luận
Bản báo cáo trước sự "chuyên minh" chậm	(ĐCSVN) – Nhiều đại biểu HĐND TP Hà Nội bày tỏ "bản báo cáo" khi các	Báo Đảng Cộng Sản	https://dangcongson.vn/thoi-su/ban-bao-cao-truoc-su-chuyen-minh-cham-chap-cua	Công bản luận
Phạt hay trước sự việc, thì phải gọi là gia đình	(ĐCSVN) – Học sinh ở nhà trường phải được thông tin trước luật gia đình	Báo Đảng Cộng Sản	https://dangcongson.vn/giao-duc/phat-hay-truoc-su-viec-phoi-hop-gia-dinh-va	Công bản luận
Nếu trẻ trai là đồng nghĩa còn vấn đề về	(ĐCSVN) – Nhiều người Hà Nội sẽ mất minh chứng từ để cùng việc không	Báo Đảng Cộng Sản	https://dangcongson.vn/tieu-diem/oi-nhieu-nguoi-ha-noi-se-mat-minh-chung-tu-de	Công bản luận
Hành thiện thì chết - Khẩu độ phá chiến	LT.S. Quyết tâm hoàn thiện thể chế trong công tác xây dựng Đảng, phòng	Báo Đảng Cộng Sản	https://dangcongson.vn/thoi-su/hanh-thien-ke-chet-khou-do-pha-chien-lua	Công bản luận
Giao pháp nào cho các chuyên "biết" rồi	(ĐCSVN) – Nhiều nhà văn, tác gia giao thông đã trở thành văn nhân, trước	Báo Đảng Cộng Sản	https://dangcongson.vn/thoi-su/giao-phap-nao-cho-cac-chuyen-biet-roi	Công bản luận
Nhân văn, nhân đạo trong kỷ nguyên số	(ĐCSVN) – Tổng Bí thư Nguyễn Phú Trọng nhiều lần khẳng định các vị	Báo Đảng Cộng Sản	https://dangcongson.vn/thoi-su/nhan-van-nhan-dao-trong-ky-nguyen-so	Công bản luận
Phong cách làm việc "Xây" văn đề pháp	(ĐCSVN) – Là người phải nghĩ trước các xã hội không thể nghĩ những	Báo Đảng Cộng Sản	https://dangcongson.vn/thoi-su/phong-cach-lam-viec-xay-van-de-phap	Công bản luận
Hội nghị COP27: Kyong won về chất thải	(ĐCSVN) – Hội nghị lần thứ 27 các bên tham gia Công ước khung của Liê	Báo Đảng Cộng Sản	https://dangcongson.vn/giao-duc/hoi-nghi-cop27-kyong-won-ve-chat-thai	Công bản luận
Xóa bỏ sản xuất mìn mìn, nhỏ lệ	(ĐCSVN) – Sản xuất mìn mìn, nhỏ lệ là vấn đề lớn mà hội nghị đến nay	Báo Đảng Cộng Sản	https://dangcongson.vn/kinh-te/xoa-bo-san-xuat-min-min-nho-le-626035.htm	Công bản luận
Là "nóng" xuyên các đống đá của World	(ĐCSVN) – Liên tiếp trong những ngày qua, Bộ Công an, cơ quan công	Báo Đảng Cộng Sản	https://dangcongson.vn/cung-ban-luan/la-nong-xuyen-cac-dong-da-cua-world	Công bản luận
Đảm bảo an toàn thực phẩm trong các	(ĐCSVN) – Những ngày gần đây, thông tin về hội sinh tại trường (School	Báo Đảng Cộng Sản	https://dangcongson.vn/cung-ban-luan/dam-bao-an-toan-thuc-pham-trong-cac-b	Công bản luận
Quyền lực của khán giả	(ĐCSVN) – Khán giả là hiện tượng lực của minh chứng qua việc từ 0 thời	Báo Đảng Cộng Sản	https://dangcongson.vn/cung-ban-luan/quyen-luc-cua-khan-gia-625770.htm	Công bản luận
Gỡ bỏ thuế giảm áp lực đã hạn tại	(ĐCSVN) – Luật lệ trong 9 tháng đầu năm 2022, tổng giá trị tại trước đã	Báo Đảng Cộng Sản	https://dangcongson.vn/cung-ban-luan/goi-bo-thue-giam-ap-luc-da-han-trai-phie	Công bản luận
TP Hồ Chí Minh đặt minh chuyển đổi số	(ĐCSVN) – Những năm qua, quy mô ngành giao vận và đạo tạo của TP H	Báo Đảng Cộng Sản	https://dangcongson.vn/cung-ban-luan/tp-ho-chi-minh-da-phat-minh-chuyen-doi-so	Công bản luận
Vị sự an toàn trên một con đường	(ĐCSVN) – Tại nạn giao thông có thể phòng tránh được. Để không bị	Báo Đảng Cộng Sản	https://dangcongson.vn/tieu-diem/vi-su-an-toan-trien-mot-con-duong-624918.htm	Công bản luận
Cơ hội lịch của tiêu dùng cuối năm	(ĐCSVN) – Được thực hiện trên phạm vi toàn quốc từ ngày 15/11 - 22/12	Báo Đảng Cộng Sản	https://dangcongson.vn/cung-ban-luan/co-hoi-lich-cua-tieu-dung-cuoi-nam-624687	Công bản luận
Lời nhắn nhủ của "Tư lệnh" Ngành giáo	(ĐCSVN) – Liên quan đến việc tăng lương, phụ cấp cho một công viên yên	Báo Đảng Cộng Sản	https://dangcongson.vn/tieu-diem/loi-nhan-nhu-cua-tu-lenh-nganh-giao-duc	Công bản luận
"Sáng lọc" các dự án FDI tiềm ẩn rủi	(ĐCSVN) – Dù dòng vốn đầu tư nước ngoài (FDI) là một trong các nguồ	Báo Đảng Cộng Sản	https://dangcongson.vn/cung-ban-luan/sang-loc-cac-du-an-fdi-tiem-an-rui-ro-6247	Công bản luận
Kiểm tra, thực hiện, sáng suốt	(ĐCSVN) – Sau thời gian 21 ngày làm việc, Kỳ họp thứ tư Quốc hội khóa X	Báo Đảng Cộng Sản	https://dangcongson.vn/ban-diem/kiem-tra-thuc-hien-sang-suot	Công bản luận
Tâm lực, trí lực trong công nghệ quốc	(ĐCSVN) – Kịp thời thực thi, Quốc hội khóa XV đã vượt qua những	Báo Đảng Cộng Sản	https://dangcongson.vn/ban-diem/tam-luc-tri-luc-trong-cong-nghe-quoc	Công bản luận
Cả hai trách nhiệm người đứng đầu	(ĐCSVN) – Cả hai trách nhiệm người đứng đầu trong giai đoạn	Báo Đảng Cộng Sản	https://dangcongson.vn/cung-ban-luan/ca-hai-trach-nhiem-nguoi-dung-dau-tr	Công bản luận
Thu phí thứ, chưa đủ	(ĐCSVN) – Thông tin về việc TP Hà Nội sẽ lấy tiền kinh phí và việc thu	Báo Đảng Cộng Sản	https://dangcongson.vn/cung-ban-luan/thu-phi-thu-thu-cho-da-624502.htm	Công bản luận
Đầu tư thực, đúng chỉ chấm cần vào	(ĐCSVN) – Các đại biểu Quốc hội cho rằng, nền tảng xây dựng	Báo Đảng Cộng Sản	https://dangcongson.vn/vu-to/dao-tu-thuc-dung-chi-cham-can-va-loai-nao	Công bản luận
Số cơ sở tham tra, giám sát kỹ lưỡng về	(ĐCSVN) – Một dự án cơ sở cơ sở tham tra giám sát kỹ lưỡng đã, song	Báo Đảng Cộng Sản	https://dangcongson.vn/ban-diem/su-co-so-tham-tra-giam-sat-ky-luong-ve	Công bản luận
Không ngừng với thông qua luật, không	(ĐCSVN) – Không ngừng với thay thế luật, không "hành chính" mà	Báo Đảng Cộng Sản	https://dangcongson.vn/tieu-diem/khong-nguoi-thong-qua-luat-thanh-lam-luat-c	Công bản luận
Quy định rõ điều kiện thu hồi đất, tranh	(ĐCSVN) – Thúc tác cho thực, văn đề thu hồi đất làm vấn đề phức	Báo Đảng Cộng Sản	https://dangcongson.vn/thoi-su/auy-dinh-ro-dieu-kien-thu-hoi-dat-tranh-an-dung-c	Công bản luận

Hình 2: Một số bài báo trong tập ngữ liệu



3 Giới thiệu phương pháp chính

Thay vì sử dụng một mô hình đơn lẻ để xếp hạng tất cả các tài liệu, phương pháp **Multi-stage Document Ranking** áp dụng một chuỗi các mô hình khác nhau, mỗi mô hình thực hiện một chức năng cụ thể.

3.1 Quy trình xử lý của Multi-stage Document Ranking

Đầu tiên, ta sử dụng một **phương pháp term-based** như TF-IDF hay BM25 để lọc các bài báo có nội dung không liên quan.

Ở đây, TF-IDF được sử dụng với mục đích chuyển đổi các query cùng với corpus (nội dung bài báo bao gồm title và abstrac) về dạng feature vector. Tiếp đó, các bài báo này sẽ được xếp hạng dựa trên độ đo **cosine similarity**.

Khác với TF-IDF, BM25 không tính toán feature vector. Thay vào đó, thuật toán trực tiếp xếp hạng các bài báo dựa trên **BM25 score**.

Sau khi xếp hạng, ta chọn ra 50 bài báo có xếp hạng cao nhất và những bài báo này sẽ được tiếp tục đưa qua các pre-trained language model để thực hiện **reranking** - sắp xếp lại theo mức độ liên quan với câu truy vấn, cả về mặt ngữ cảnh và ngữ nghĩa.

Đối với PhoBERT, vector ẩn của token [CLS] ở đầu ra của lớp transformer cuối cùng được đưa qua một lớp dense (fully connected layer) để giảm chiều, tạo ra vector **pooler output** và ta sẽ sử dụng vector này sẽ được sử dụng như một feature vector.

Tương tự như PhoBERT, trong Vietnamese-SBERT, ta sẽ sử dụng vector tương ứng với token [CLS] ở đầu ra của lớp transformer cuối cùng làm feature vector.

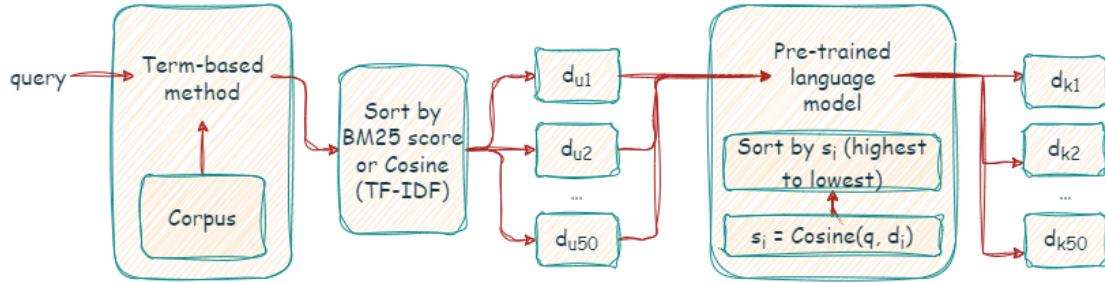
Sau khi có biểu diễn feature vector của câu query và 50 bài báo trước đó, quá trình reranking sẽ được thực hiện, sử dụng độ đo **cosine similarity**.

Về vai trò, các phương pháp term-based TF-IDF hay BM25 được sử dụng trước để lọc bớt những bài báo có nội dung không liên quan đến truy vấn, rút ngắn thời gian xử lý cho các pre-trained language model. Tiếp đó, pre-trained language model được sử dụng kết hợp nhằm mục đích khắc phục nhược điểm không nắm bắt được ngữ nghĩa văn bản của các phương pháp term-based. Như



vậy, Multi-stage Document Ranking kết hợp, phát huy ưu điểm và khắc phục nhược điểm của từng phương pháp.

Tổng kết lại, ngữ liệu và truy vấn sau khi được tiền xử lý dữ liệu sẽ đi qua pipeline như hình dưới đây.



Hình 3: Pipeline xử lý bài toán đặt ra

3.2 Các phương pháp term-based

3.2.1 TF-IDF

TF-IDF (term frequency-inverse document frequency) là độ đo để đánh giá mức độ quan trọng của một từ trong một văn bản. TF-IDF bao gồm 2 thành phần: **TF** - Term Frequency và **IDF** - Inverse Document Frequency.

Cụ thể, với từ khóa t , tài liệu d , ta có:

$$\text{TF}(t, d) = \frac{\text{Số lần từ } t \text{ xuất hiện trong tài liệu } d}{\text{Tổng số từ trong tài liệu } d}$$

Giá trị TF đo lường tần suất xuất hiện của 1 từ trong 1 văn bản. Nếu giá trị TF càng cao thì mức độ liên quan càng cao.

$$\text{IDF}(t) = 1 + \log \left(\frac{\text{Tổng số lượng tài liệu } N}{1 + \text{Số tài liệu chứa từ } t} \right)$$

Trong công thức của IDF, 1 được thêm vào mẫu để tránh lỗi chia cho 0. Còn giá trị 1 ở ngoài là để điều chỉnh cho giá trị IDF của một từ vẫn bằng 1, kể cả khi từ đó xuất hiện trong toàn bộ tập tài liệu, tránh việc giảm quá nhiều mức độ quan trọng của từ đó.

Giá trị IDF phản ánh mức độ phổ biến của từ trong toàn bộ tập tài liệu. Nếu



một từ xuất hiện trong nhiều tài liệu (phổ biến) thì giá trị IDF của nó sẽ thấp và ngược lại.

Kết hợp TF và IDF, ta sẽ xác định được mức độ quan trọng thực sự của từ trong toàn bộ tập tài liệu:

$$\text{TF-IDF}(t, d) = \text{TF}(t, d) \times \text{IDF}(t)$$

Thực chất, việc dùng TF-IDF là để biểu diễn từng văn bản và truy vấn dưới dạng một vector, sau đó mới sử dụng một phương thức như L2 hay Cosine similarity để tính toán mức độ tương đồng. Trong đề án, **cosine similarity** là độ đo được sử dụng.

3.2.2 BM25

BM25 là một hàm xếp hạng giúp xếp hạng một tập các tài liệu dựa vào việc tính xác suất xuất hiện của các từ khóa trong câu truy vấn. BM25 được cải thiện dựa trên nền tảng của TF-IDF, xem xét thêm một số yếu tố như độ dài tài liệu, các tham số điều chỉnh, ...

Với một truy vấn q chứa các từ khóa q_1, q_2, \dots, q_n , BM25 score của một tài liệu D được xác định bởi:

$$\text{BM25}(q, D) = \sum_{i=1}^n \text{IDF}(q_i) \frac{\text{TF}(q_i, D)(k_1 + 1)}{\text{TF}(q_i, D) + k_1(1 - b + b \frac{|D|}{\text{avgdl}})}$$

- $f(q_i, D)$ là tần số q_i xuất hiện trong tài liệu D .
- $|D|$ là số lượng các từ trong tài liệu D .
- avgdl là độ dài trung bình số lượng từ của các tài liệu trong ngữ liệu.
- k_1 và b là các tham số điều chỉnh. Thông thường, k_1 nằm trong khoảng $[1.2; 2.5]$ và b được gán bằng 0.75.

Các bài báo có BM25 score cao có nghĩa là chúng chứa những từ khóa có liên quan đến truy vấn. Ngược lại, các bài báo có BM25 score thấp chứa nội dung không liên quan đến truy vấn và sẽ không được xử lý trong pre-trained language model.

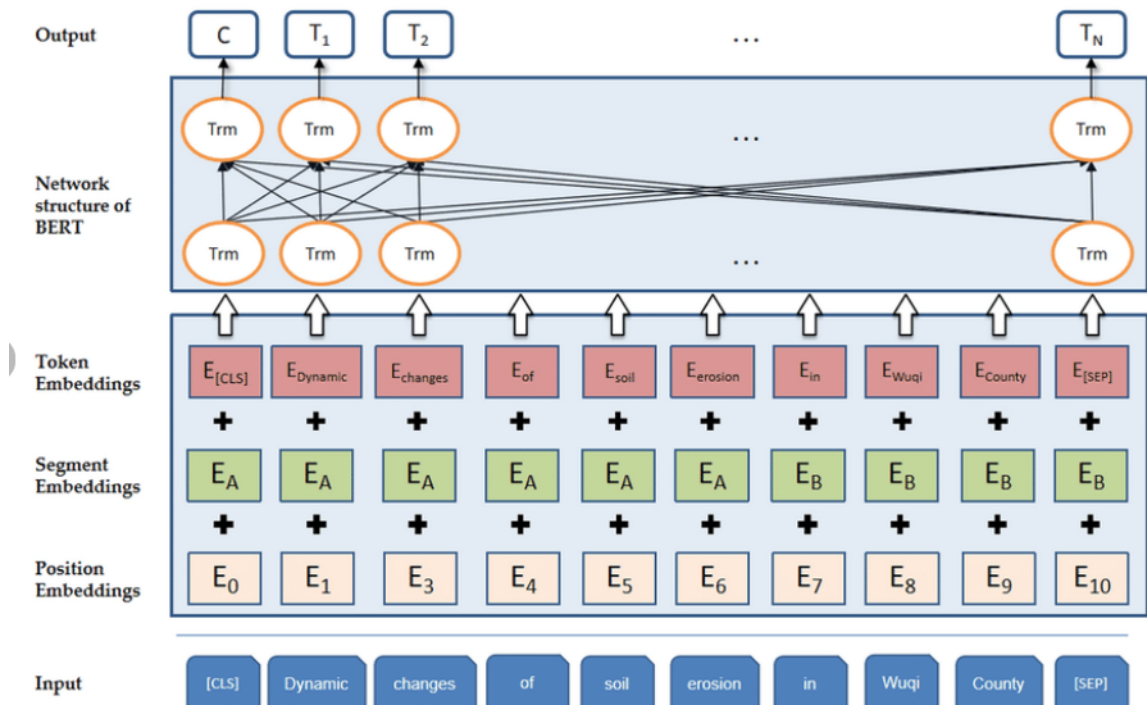


3.3 Pre-trained language model

3.3.1 PhoBERT

PhoBERT là một pre-trained language model được phát triển dựa trên BERT, nhưng được fine tune và tối ưu hóa cho tiếng Việt.

Ta sẽ sử dụng `pooled_output` của token `[CLS]` như feature vector cho câu truy vấn và nội dung của các bài báo.

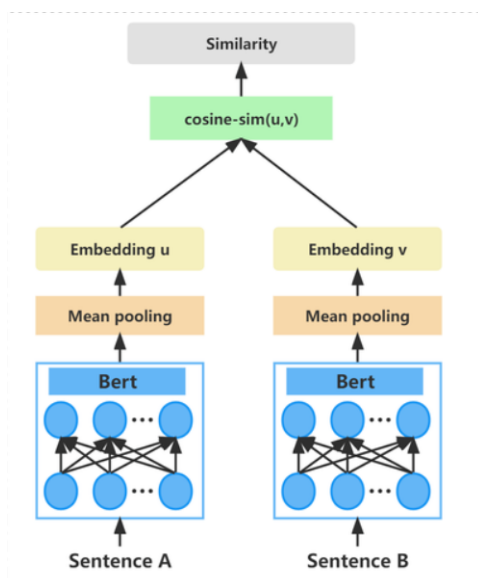


Hình 4: Cấu trúc tổng quát của BERT

3.3.2 Vietnamese-SBERT

Vietnamese-SBERT là một pre-trained language model được sử dụng để biểu diễn văn bản tiếng Việt dưới dạng feature vector, được xây dựng dựa trên mô hình SentenceBERT, cũng là một mô hình biến thể được phát triển dựa trên BERT.

Mục đích chính của mô hình này là để nhận ra được các văn bản đồng nghĩa nhau (nội dung bài báo nào gần nghĩa nhất so với câu truy vấn) thông qua các feature vector (minh họa bởi hình 5).



Hình 5: Quy trình cơ bản của SentenceBERT

4 Cài đặt phương pháp

4.1 Tiền xử lý dữ liệu

Với từng bài báo trong ngữ liệu, dữ liệu dạng text tương ứng với bài báo đó sẽ là kết hợp của title và abstract. Các truy vấn cũng như dữ liệu của các bài báo sẽ đi qua 4 bước tiền xử lý dữ liệu, bao gồm:

- Lowercasing (Di sản → di sản)
- Word segmentation (di sản → di_sản)
- Loại bỏ stopwords
- Loại bỏ các ký tự đặc biệt (tp.hcm → tphcm)

4.2 Đánh giá bài báo bằng Cosine Similarity

Để đánh giá và so sánh độ tương đồng giữa truy vấn và nội dung bài báo, nhóm sử dụng phép đo **cosine similarity** cho tác vụ này. Cụ thể như sau:



- Xét 2 vector u và v là các [CLS] đã được padding về kích thước 768x1 sau khi được xử lý qua multi-stage document ranking, độ tương đồng giữa q và u được tính bằng cách tính cosine của 2 vector này.

$$\text{Cosine Similarity}(u, v) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}$$

- Trong đó: $\mathbf{u} \cdot \mathbf{v}$ là tích vô hướng của u và v , $\|\mathbf{u}\|$ và $\|\mathbf{v}\|$ lần lượt là chuẩn 2 của vector u và v . Nếu giá trị CosSimilarity của 2 vector tương ứng với 2 văn bản càng gần 1 thì 2 văn bản này càng đồng nghĩa với nhau và ngược lại, nếu càng gần -1 thì 2 văn bản này càng có ý nghĩa ngược nhau.

5 Đánh giá phương pháp

Thực hiện so sánh độ hiệu quả của hướng tiếp cận Multi-stage document ranking cho truy vấn các bài báo tiếng Việt bằng cách kết hợp một trong các phương pháp term-based (BM25, TF-IDF) cùng với một trong các pre-trained language model (PhoBERT, Vietnamese-SBERT). So sánh tất cả các cặp có thể tạo ra, ghi lại hiệu quả và tốc độ xử lý của từng phương pháp.

Với truy vấn q , ta gọi D_{pred} là danh sách các bài báo đã được sắp xếp từ trên xuống dưới là liên quan với q sau khi sử dụng multi-stage document ranking, D_{truth} là danh sách các bài báo đã được sắp xếp từ trên xuống dưới mà ta thật sự mong muốn. Ta so sánh kết quả giữa mô hình và thực tế dựa trên 10 bài báo đầu tiên bằng $nDCG_{10}$.

Công thức của $nDCG_{10}$ được tính như sau:

$$nDCG_{10} = \frac{DCG_{10}}{IDCG_{10}}$$

trong đó:

- $DCG_{10} = \sum_{i=1}^{10} \frac{rel_i}{\log_2(i+1)}$, rel_i là độ liên quan của bài báo thứ i trong D_{pred} .
- $IDCG_{10} = \sum_{i=1}^{10} \frac{rel_i}{\log_2(i+1)}$, rel_i là độ liên quan của bài báo thứ i trong D_{truth} .
- Giá trị $nDCG_{10}$ nằm trong đoạn $[0; 1]$, càng gần 1 thì mô hình càng tốt.



5.1 Tập đánh giá

Tập đánh giá gồm 300 bài báo đã được xếp hạng và gán nhãn mức độ liên quan (dành cho việc tính $nDCG_{10}$) với truy vấn: "Giá lăn bánh của Vinfast Vf8 năm 2024?"

title	abstract	link	rank	score
Giá xe VinFast VF 8 mới nhất tháng 7/2024	Cập nhật giá xe VinFast VF 8 2024 kèm tin khuyến mại, thông	https://oto.com.vn/bang-gia-xe-o-to-vinfast-vf-8-moi-nhat-thang-7-2024	1	4
Bảng giá xe VinFast kèm ưu đãi mới nhất tháng 7/2024	Bảng giá xe VinFast 2024 kèm tin khuyến mại, hình ảnh, video	https://oto.com.vn/bang-gia-xe-o-to-vinfast-moi-nhat-thang-7-2024	2	4
Bảng giá xe VinFast tháng 4/2023 kèm ưu đãi	Cập nhật bảng giá xe VinFast mới nhất tháng 4/2023 tại Việt	https://vinfast.vn/bang-gia-xe-vinfast-thang-4-2023-ken	3	4
Vinfast VF8 2024: Giá xe Vinfast VF8 mới nhất và những thông tin	Được trang bị nhiều công nghệ hiện đại, Vinfast VF8 hứa hẹn	https://tinxe.vn/gia-xe-vinfast-vf8	4	4
Giá xe VinFast VF 8 tháng 11/2023: TSKT & đánh giá chi tiết	Tham khảo ngay giá lăn bánh mới nhất của mẫu xe Vinfast VF	https://thethao247.vn/364-gia-xe-vinfast-vf-8-d250055	5	4
Xe Điện VinFast VF8 2024: Bảng Giá Lăn Bánh, Thông số kĩ thuật	Cùng VinFast Sài Gòn cập nhập các thông tin mới nhất về xe	https://oto-vinfastsaigon.com/xe-vinfast-vf8/7gad_sour	6	4
So sánh Hyundai Ioniq 5 hay VinFast VF8 tầm giá 1,5 tỷ đồng	Hai mẫu xe điện lắp ráp tại thị trường Việt Nam đang được nh	https://www.24h.com.vn/o-to/so-sanh-hyundai-ioniq-5	7	4
Sau 1 năm lăn bánh, VinFast VF8 được rao bán khá khó tin	Những chiếc VinFast VF8 đang được rao bán khá nhiều hiện	https://danviet.vn/sau-1-nam-lan-banh-vinfast-vf8-duo	8	3
Xe điện VinFast VF 8 lăn bánh hơn 1 năm "hết giá" chỉ chưa đến 80	Ở thời điểm hiện tại, những chiếc xe điện đã qua sử dụng của	https://auto5.vn/325-xe-dien-vinfast-vf-8-lan-banh-hon	9	3
SUV điện VinFast VF 8 rớt giá bao nhiêu sau 2 năm lăn bánh?	Một chiếc SUV điện VinFast VF 8 hiện đang được một cơ sở k	https://thethao247.vn/450-suv-dien-vinfast-vf-8-rot-gia	10	3
VinFast VF8 Plus lần đầu xuống giá chỉ hơn 800 triệu đồng sau 1 n	VinFast VF8 Plus là bản cao cấp từng có giá hơn 1 tỷ đồng trê	https://etime.danviet.vn/vinfast-vf8-plus-lan-dau-xuong	11	3
VinFast VF8 mất giá khó tin sau 1 năm lăn bánh	Sau 1 năm lăn bánh, VinFast VF8 hiện đang được rao bán kh	https://auto5.vn/326-vinfast-vf8-mat-gia-kho-tin-sau-1	12	3
Sau 2 năm lăn bánh, VinFast VF 8 "lướt" lên sàn xe cũ với mức giá	Theo tham khảo, một chiếc VinFast VF 8 sản xuất năm 2022	https://thethao247.vn/419-sau-2-nam-lan-banh-vinfast	13	3

Hình 6: Dữ liệu tập đánh giá

Quy tắc gán nhãn mức độ liên quan được quy định như sau:

- 0: Bài báo không chứa nội dung liên quan đến truy vấn.
- 1: Bài báo chứa thông tin giúp ta có được các nội dung bên lề (một phần các chủ đề) liên quan đến truy vấn.
- 2: Bài báo chứa thông tin hữu ích giúp ta có được các nội dung chứa khía cạnh liên quan đến các chủ đề của truy vấn.
- 3: Bài báo chứa các thông tin mang tính cụ thể hơn so với mức 2 (tức là hơn cả về mặt chủ đề, nó mang đến các thông tin quan trọng về người, sự kiện nào đó...), sao cho giải thích được một phần truy vấn.
- 4: Bài báo chứa các ngữ cảnh quan trọng chứa đựng hoặc giải thích được đầy đủ các thông tin liên quan đến truy vấn.

Mức độ liên quan giữa bài báo ứng với truy vấn "Giá lăn bánh của vinfast vf8 năm 2024" được thể hiện trên cột "score" của tập đánh giá. Trong tập đánh giá, có 7 bài báo có relevance score là 4, 8 bài báo là 3, 7 bài báo là 2, 15 bài báo là 1 và 13 bài còn lại là 0.



5.2 Kết quả thu được

Giá trị $nDCG_{10}$ của từng hướng tiếp cận trong bảng dưới (kết quả được làm tròn đến chữ số thập phân thứ 4). Kết quả cao nhất của bảng được **in đậm**. Các kết quả được lấy trung bình trên **10 lần chạy** và **không sử dụng GPU**.

Bảng 1: Kết quả $nDCG_{10}$ của một số phương pháp

Hướng tiếp cận	Phương pháp	$nDCG_{10}$	Thời gian thực thi
Term-based	TF-IDF	0.9090	0.02 s
	BM25	0.8541	0.04 s
Pre-trained language model	Vietnamese-SBERT	0.8399	58.85 s
	PhoBERT	0.9170	60.91 s
Multi-stage document ranking	TF-IDF + Vietnamese-SBERT	0.8399	10.34 s
	TF-IDF + PhoBERT	0.9513	11.43 s
	BM25 + Vietnamese-SBERT	0.8002	11.86 s
	BM25 + PhoBERT	0.8268	11.33 s

Nhận xét:

- Hướng tiếp cận **TF-IDF kết hợp PhoBERT** cho kết quả tốt nhất với điểm $nDCG_{10}$ đạt đến **0.9513**
- Có thể thấy khi sử dụng multi-stage document ranking thay vì chỉ sử dụng đơn lẻ pre-trained language model là PhoBERT thì kết quả và tốc độ xử lý đều tốt hơn.
- Phương pháp **TF-IDF** cho điểm $nDCG_{10}$ khá tốt với giá trị **0.91**, đặc biệt hiệu quả khi tốc độ thực thi rất nhanh (chỉ mất **0.02s**).
- Từ đây, ta có nhận xét: khi cần ưu tiên tính real-time (tức là tốc độ trả về output truy vấn phải gần như ngang bằng tốc độ đưa vào input), ta có thể đánh đổi độ chính xác như phương pháp TF-IDF. Ngược lại, nếu cần ưu tiên độ chính xác của kết quả truy vấn, việc kết hợp các pre-trained language model sẽ cho kết quả tối ưu.



6 Tài liệu tham khảo

1. PhoBERT: Pre-trained language models for Vietnamese - <https://aclanthology.org/2020.findings-emnlp.92.pdf>
2. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding - <https://arxiv.org/pdf/1810.04805.pdf>
3. Improvements to BM25 and Language Models Examined - <https://www.cs.otago.ac.nz/homepages/andrew/papers/2014-2.pdf>
4. Rank-BM25: A two line search engine - <https://pypi.org/project/rank-bm25/>
5. IR-BERT: Leveraging BERT for Semantic Search in Background Linking for News Articles - <https://arxiv.org/abs/2007.12603>
6. Ad-hoc retrieval with BERT - <https://arxiv.org/abs/2007.12603>
7. Vietnamese-SBERT: <https://huggingface.co/keepitreal/vietnamese-sbert>
8. TF-IDF là gì?: <https://vi.wikipedia.org/wiki/Tf%E2%80%93idf>
9. Thực thi TF-IDF: https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.TfidfVectorizer.html