

TRƯỜNG ĐẠI HỌC SƯ PHẠM KỸ THUẬT TP. HỒ CHÍ MINH  
KHOA CÔNG NGHỆ THÔNG TIN



**TIỂU LUẬN CUỐI KỲ**

**Môn học: Xử lý ngôn ngữ tự nhiên**

**Tên tiểu luận: Xây dựng Chatbot tích hợp hệ khuyến nghị Item-based và nhận diện cảm xúc**

Giảng viên: TS. Phan Thị Huyền Trang

**Danh sách sinh viên thực hiện**

Mã số SV	Họ và tên	Mức độ đóng góp (%)
22110138	Nguyễn Trung Hiếu	100%
22110200	Lê Hoàng Bảo Phúc	100%

*TP. Hồ Chí Minh, tháng 5 năm 2025*

# MỤC LỤC

DANH MỤC TỪ VIẾT TẮT.....	1
DANH MỤC HÌNH ẢNH .....	2
DANH MỤC BẢNG BIỂU .....	3
PHẦN 1: MỞ ĐẦU.....	4
1.1. Phát biểu đề tài.....	4
1.2. Mục đích, yêu cầu cần thực hiện.....	5
1.3. Phạm vi và đối tượng .....	6
PHẦN 2: CƠ SỞ LÝ THUYẾT.....	8
2.1. Chatbot .....	8
2.1.1. Khái niệm Chatbot .....	8
2.1.2. Lợi ích và ứng dụng của Chatbot.....	8
2.1.3. Các loại Chatbot.....	9
2.2. Item-based Recommendation System.....	10
2.2.1. Khái niệm .....	10
2.2.2. Cosine Similarity.....	11
2.3. Logistic Regression (LR) .....	12
2.3.1. Khái niệm .....	12
2.3.2. Multinomial Logistic Regression.....	17
PHẦN 3: PHÂN TÍCH, THIẾT KẾ.....	19
3.1. Ý tưởng thiết kế.....	19
3.2. Công nghệ được áp dụng .....	20
3.2.1. Django .....	20
3.2.2. VADER .....	21
3.3. Tài nguyên được sử dụng.....	23
3.3.1. Tập dữ liệu .....	23
3.3.1.1. Bitext Gen AI Chatbot Customer Support Dataset .....	23
3.3.1.2. Consumer Reviews of Amazon Products.....	23
3.3.2. Thư mục .....	25
3.3.2.1. Thư mục Dataset .....	25
3.3.2.2. Thư mục Sentiment .....	26
3.3.2.3. Thư mục Collaborative Recommendation .....	27
3.3.2.4. Thư mục Sentiment_User .....	28

3.3.2.5. Thư mục Web.....	28
PHẦN 4: THỰC NGHIỆM.....	30
4.1. Kết quả thực nghiệm .....	30
4.2. Giao diện .....	34
PHẦN 5: KẾT LUẬN.....	37
5.1. Kết luận .....	37
5.2. Hướng phát triển .....	38
TÀI LIỆU THAM KHẢO.....	39

## DANH MỤC TỪ VIẾT TẮT

Thuật ngữ	Ý nghĩa
LR	Logistic Regression

## DANH MỤC HÌNH ẢNH

Hình 4.1.1: Thư Mục Dataset .....	30
Hình 4.1.2: Intent Và Category .....	30
Hình 4.1.3: Thư Mục Sentiment.....	31
Hình 4.1.4: Thư Mục Collaborative Recommendation .....	32
Hình 4.1.5: Thư Mục Sentiment_User .....	32
Hình 4.1.6: Báo Cáo Phân Loại.....	32
Hình 4.1.7: Accuracy Của Cột Training Và Cột Test .....	33
Hình 4.1.8: Ma Trận Nhầm Lẫn .....	34
Hình 4.2.1: Giao Diện Bắt Đầu .....	34
Hình 4.2.2: Tư Vấn Và Gợi Ý Sản Phẩm.....	35
Hình 4.2.3: Nhận Diện Cảm Xúc Tích Cực .....	35
Hình 4.2.4: Trang Chủ Web .....	36

## **DANH MỤC BẢNG BIỂU**

Bảng 3.3.1.1.1: Danh Sách Cột Và Thuộc Tính.....	23
Bảng 3.3.1.2.1: Danh Sách Cột Và Ý Nghĩa.....	24
Bảng 3.3.2.3.1: Hai Phương Án Xây Dựng Hệ Khuyến Nghị .....	27

# PHẦN 1: MỞ ĐẦU

## 1.1. Phát biểu đề tài

Trong bối cảnh công nghệ số phát triển vượt bậc, các hệ thống giao tiếp tự động đang ngày càng khẳng định vai trò quan trọng trong việc kết nối doanh nghiệp với khách hàng. Những công cụ này không chỉ cần đảm bảo khả năng cung cấp thông tin nhanh chóng mà còn phải đáp ứng kỳ vọng về sự thấu hiểu và cá nhân hóa trong tương tác. Nhu cầu xây dựng một hệ thống có khả năng nhận biết cảm xúc của người dùng thông qua văn bản và đưa ra các đề xuất phù hợp với sở thích cá nhân đã trở thành một hướng đi đầy tiềm năng. Dự án được thực hiện nhằm khám phá việc tích hợp các tính năng thông minh này, hướng đến việc tạo ra một công cụ giao tiếp tiên tiến, đáp ứng tốt hơn nhu cầu của người dùng trong các lĩnh vực như thương mại điện tử, dịch vụ khách hàng, và hỗ trợ cá nhân.

Về mặt khoa học, việc nghiên cứu một hệ thống giao tiếp tự động tích hợp khả năng phân tích cảm xúc và khuyến nghị cá nhân hóa mang lại nhiều giá trị đáng kể. Nghiên cứu này góp phần làm rõ cách các kỹ thuật trí tuệ nhân tạo, đặc biệt là xử lý ngôn ngữ tự nhiên và học máy, có thể được kết hợp để tạo ra một hệ thống giao tiếp thông minh hơn. Việc khám phá sự tương tác giữa các thành phần như nhận diện trạng thái tâm lý và đề xuất sản phẩm dựa trên hành vi người dùng không chỉ mở rộng hiểu biết về công nghệ chatbot mà còn đặt nền tảng cho các nghiên cứu tiếp theo trong lĩnh vực giao tiếp người-máy. Hơn nữa, dự án cung cấp cơ hội để đánh giá hiệu quả của việc tích hợp các phương pháp phân tích dữ liệu phức tạp vào một hệ thống thống nhất, từ đó đóng góp vào sự phát triển của các giải pháp công nghệ tiên tiến.

Về mặt thực tiễn, hệ thống được đề xuất hứa hẹn mang lại nhiều lợi ích cho cả doanh nghiệp và người dùng. Trong môi trường kinh doanh cạnh tranh, việc cung cấp các tương tác được cá nhân hóa dựa trên cảm xúc và sở thích của khách hàng có thể cải thiện đáng kể trải nghiệm người dùng, từ đó tăng cường sự hài lòng và lòng trung thành. Chẳng hạn, trong lĩnh vực thương mại điện tử, một hệ thống có khả năng nhận biết khi khách hàng đang bối rối hoặc không hài lòng và đưa ra các gợi ý phù hợp có thể nâng cao chất lượng dịch vụ và thúc đẩy doanh thu. Đối với dịch vụ khách hàng, công cụ này có thể giúp xử lý các yêu cầu một cách nhạy bén hơn, giảm thời gian phản hồi và cải thiện hiệu

quả vận hành. Ngoài ra, tính ứng dụng của hệ thống không chỉ giới hạn ở lĩnh vực kinh doanh mà còn có thể mở rộng sang các ngành như giáo dục, y tế, hoặc hỗ trợ cá nhân, nơi mà sự thấu hiểu và phản hồi phù hợp đóng vai trò quan trọng.

Dự án tập trung vào việc khai thác tiềm năng của các công nghệ hiện đại để đáp ứng nhu cầu ngày càng cao của người dùng trong thời đại số. Bằng cách tích hợp khả năng phân tích cảm xúc với cơ chế khuyến nghị dựa trên dữ liệu thực tế, hệ thống hướng đến việc tạo ra một công cụ giao tiếp không chỉ thông minh mà còn mang tính nhân văn, có thể hiểu và đáp ứng tốt hơn nhu cầu của từng cá nhân. Việc phát triển một hệ thống như vậy không chỉ phản ánh xu hướng công nghệ hiện nay mà còn thể hiện tầm quan trọng của việc kết hợp giữa khoa học dữ liệu và trải nghiệm người dùng trong việc định hình tương lai của các dịch vụ số.

## **1.2. Mục đích, yêu cầu cần thực hiện**

Dự án này sẽ tạo ra một hệ thống giao tiếp tự động có khả năng phân tích tâm trạng người dùng qua văn bản và đề xuất các sản phẩm hoặc dịch vụ phù hợp với sở thích cá nhân, hướng đến nâng cao trải nghiệm người dùng và tối ưu hóa hiệu quả hoạt động trong các lĩnh vực như thương mại điện tử và dịch vụ khách hàng. Hệ thống được thiết kế để đáp ứng nhu cầu về tương tác cá nhân hóa, giúp người dùng cảm thấy được thấu hiểu thông qua các phản hồi nhạy bén và nhận được gợi ý chính xác, đồng thời hỗ trợ doanh nghiệp cải thiện chất lượng dịch vụ, tăng sự hài lòng của khách hàng và giảm chi phí vận hành.

Hệ thống cần đáp ứng các yêu cầu cụ thể để đảm bảo chất lượng và tính thực tiễn. Nó phải phân tích nội dung giao tiếp để xác định chính xác trạng thái tâm lý của người dùng, như vui, buồn hay trung lập, và điều chỉnh phản hồi sao cho phù hợp, tạo cảm giác tự nhiên. Hệ thống cũng cần đưa ra các gợi ý sản phẩm hoặc dịch vụ dựa trên dữ liệu hành vi, đảm bảo tính phù hợp với nhu cầu riêng của từng người, ví dụ như đề xuất sản phẩm trong quá trình mua sắm trực tuyến. Về hiệu suất, hệ thống phải xử lý nhanh, duy trì độ chính xác cao khi làm việc với khối lượng dữ liệu lớn hoặc nội dung phức tạp, đảm bảo thời gian phản hồi tối ưu để không làm gián đoạn trải nghiệm người dùng. Hơn nữa, công cụ cần có giao diện trực quan, dễ tích hợp vào các nền tảng như website hoặc ứng dụng di động, phù hợp với các tình huống thực tế như tư vấn khách hàng hoặc hỗ trợ



bán hàng. Để đảm bảo tính bền vững, hệ thống phải được thiết kế linh hoạt, cho phép cập nhật dữ liệu mới hoặc mở rộng ứng dụng sang các lĩnh vực khác như giáo dục hoặc chăm sóc sức khỏe mà không cần thay đổi cấu trúc cốt lõi. Những yêu cầu này đảm bảo hệ thống không chỉ hoạt động hiệu quả về mặt kỹ thuật mà còn mang lại giá trị thực tiễn, đáp ứng nhu cầu cá nhân hóa và nâng cao chất lượng dịch vụ trong môi trường số.

### **1.3. Phạm vi và đối tượng**

Dự án tập trung vào việc thiết kế và phát triển một hệ thống chatbot thông minh tích hợp hai thành phần chính: hệ khuyến nghị dựa trên mặt hàng (Item-based Recommendation System) và khả năng nhận diện cảm xúc của người dùng thông qua văn bản. Về mặt kỹ thuật, phạm vi nghiên cứu bao gồm việc ứng dụng các công nghệ xử lý ngôn ngữ tự nhiên (NLP), học máy (Machine Learning) và phân tích dữ liệu để xây dựng một hệ thống giao tiếp tự động có khả năng hiểu và phản hồi phù hợp với trạng thái tâm lý cũng như sở thích cá nhân của người dùng. Cụ thể, hệ thống sẽ sử dụng các thuật toán như Cosine Similarity để tính toán độ tương đồng giữa các mặt hàng trong hệ khuyến nghị và các phương pháp phân tích cảm xúc như VADER hoặc Multinomial Logistic Regression để xác định trạng thái tâm lý (tích cực, tiêu cực, trung tính) từ văn bản đầu vào.

Về mặt ứng dụng, dự án giới hạn trong các lĩnh vực có nhu cầu tương tác cao với khách hàng, chủ yếu là thương mại điện tử và dịch vụ khách hàng trực tuyến. Trong thương mại điện tử, chatbot sẽ hỗ trợ người dùng trong quá trình mua sắm bằng cách đưa ra các gợi ý sản phẩm dựa trên lịch sử tương tác và điều chỉnh cách phản hồi dựa trên cảm xúc được nhận diện, ví dụ như an ủi khi khách hàng thể hiện sự không hài lòng hoặc khuyến khích khi họ tỏ ra hứng thú. Trong dịch vụ khách hàng, hệ thống sẽ xử lý các yêu cầu phổ biến như tra cứu thông tin đơn hàng, giải đáp thắc mắc hoặc hỗ trợ kỹ thuật, đồng thời đảm bảo phản hồi mang tính cá nhân hóa và nhạy bén với trạng thái tâm lý của người dùng. Ngoài ra, dự án cũng xem xét khả năng mở rộng sang các lĩnh vực khác như giáo dục (hỗ trợ học tập cá nhân hóa) hoặc y tế (tư vấn sức khỏe tâm lý), nhưng trong giai đoạn này, trọng tâm vẫn là tối ưu hóa hiệu quả trong các ứng dụng thương mại.

Về công nghệ, dự án sử dụng framework Django để xây dựng giao diện web và tích hợp các thành phần của hệ thống, đảm bảo tính bảo mật, dễ mở rộng và khả năng triển khai thực tế. Các tập dữ liệu được chọn đảm bảo tính thực tiễn và phù hợp với mục tiêu xây dựng một chatbot có khả năng xử lý các tình huống thực tế. Phạm vi nghiên cứu không bao gồm việc phát triển các mô hình ngôn ngữ lớn (LLM) từ đầu mà tập trung vào việc tinh chỉnh (fine-tuning) hoặc sử dụng các mô hình, công cụ sẵn có như VADER để tối ưu hóa hiệu suất trong các tác vụ cụ thể.

Đối tượng chính của đề tài là các doanh nghiệp hoạt động trong lĩnh vực thương mại điện tử và dịch vụ khách hàng, nơi nhu cầu về tương tác nhanh chóng, hiệu quả và cá nhân hóa với người dùng là yếu tố then chốt. Cụ thể, các doanh nghiệp vừa và nhỏ, cũng như các nền tảng bán lẻ trực tuyến lớn, là những đối tượng tiềm năng có thể hưởng lợi từ hệ thống chatbot này. Ví dụ, các công ty như Shopee, Lazada hoặc Amazon có thể tích hợp hệ thống để cải thiện trải nghiệm mua sắm, tăng tỷ lệ chuyển đổi và giảm chi phí hỗ trợ khách hàng. Ngoài ra, các tổ chức cung cấp dịch vụ trực tuyến, chẳng hạn như ngân hàng, công ty viễn thông hoặc dịch vụ công nghệ, cũng nằm trong nhóm đối tượng mục tiêu, bởi họ thường xuyên xử lý khối lượng lớn yêu cầu từ khách hàng và cần các giải pháp tự động hóa thông minh.

Về mặt học thuật, đề tài hướng đến các nhà nghiên cứu, sinh viên và chuyên gia trong lĩnh vực trí tuệ nhân tạo, xử lý ngôn ngữ tự nhiên và hệ thống khuyến nghị. Dự án cung cấp một trường hợp nghiên cứu thực tiễn về việc kết hợp các kỹ thuật phân tích cảm xúc và khuyến nghị trong một hệ thống thống nhất, từ đó đóng góp vào việc phát triển các phương pháp giao tiếp người-máy tiên tiến hơn. Các kết quả từ dự án, bao gồm phân tích hiệu suất của hệ thống và các bài học kinh nghiệm trong việc tích hợp các công nghệ, có thể được sử dụng làm tài liệu tham khảo cho các nghiên cứu tiếp theo hoặc các dự án phát triển hệ thống thông minh.

## PHẦN 2: CƠ SỞ LÝ THUYẾT

### 2.1. Chatbot

#### 2.1.1. Khái niệm Chatbot

Chatbot là chương trình có thể trò chuyện với người dùng bằng văn bản(hoặc giọng nói). Các chatbot sẽ mô phỏng cuộc trò chuyện của con người và cố gắng trả lời các câu hỏi của người dùng một cách tự nhiên và giống con người nhất, chatbot hiện nay được ứng dụng nhiều nhất trong việc hỗ trợ bán hàng và hỗ trợ dịch vụ khách hàng trực tuyến. Hiện nay chatbot đã rất phổ biến và hiện đại khi sử dụng các công nghệ xử lý ngôn ngữ tự nhiên để hiểu người dùng và trả lời các câu hỏi phức tạp và có chiều sâu. Chatbot hiện nay được xây dựng không chỉ để trả lời câu hỏi của người dùng mà còn phải có khả năng cá nhân hóa trong mọi khía cạnh.

#### 2.1.2. Lợi ích và ứng dụng của Chatbot

Chatbot có thể tìm kiếm và truy xuất thông tin từ bất kỳ nơi nào ở cơ sở dữ liệu nội bộ hoặc bên ngoài và cung cấp câu trả lời thông qua cuộc trò chuyện với con người với các lợi ích:

- **Tự động hóa:** Chatbot tiết kiệm thời gian và công sức cho tổ chức, chúng kết hợp các bước từ các quy trình để tự động hóa các tác vụ từ đây xử lý yêu cầu để giải quyết các vấn đề phổ biến và có thể điều chỉnh hoạt động của chúng khi cần thiết.
- **Linh hoạt:** Chatbot phản hồi lại từ cuộc trò chuyện với người dùng, chatbot có thể được nhúng và quy trình hoạt động của phần mềm để tương tác với khách hàng. Chatbot dịch vụ khách hàng hiện nay có thể trả lời thắc mắc của họ trên cách nền tảng truyền thông xã hội, web hoặc ứng dụng nhắn tin và thậm chí có thể thiết lập trên những ứng dụng nội bộ.
- **Tương tác khách hàng:** Dịch vụ chăm sóc khách hàng chỉ dựa vào sự tương tác của con người sẽ có năng suất hạn chế và thiếu tính linh hoạt. Với chatbot có thể dễ dàng đưa ra phản hồi nhanh chóng, cá nhân hóa tương tác với khách hàng trên quy mô lớn và có thể tiếp cận khách hàng dễ dàng hơn.

Các tổ chức trong các ngành nghề sử dụng chatbot để hợp lý hóa trải nghiệm của khách hàng, tăng hiệu quả vận hành và giảm chi phí, cụ thể:

- **Tăng năng suất doanh nghiệp:** Chúng ta có thể tích hợp chatbot với các hệ thống backend của doanh nghiệp như quản lý quan hệ khách hàng(CRM), chương trình quản lý hàng tồn kho hoặc hệ thống nhân sự(HR). Chatbot có thể kiểm tra doanh số bán hàng hoặc tình trạng tồn kho, tạo báo cáo,...
- **Trợ lý cá nhân:** Chatbot có thể đơn giản hóa và đẩy nhanh các hoạt động cá nhân hàng ngày.
- **Ứng dụng trung tâm chăm sóc khách hàng:** Chatbot có thể giải quyết các yêu cầu khách hàng trong việc chăm sóc khách hàng, nhiều tác vụ có thể được chatbot giải quyết ví dụ: Lên lịch, thay đổi mật khẩu,... Các bot có thể linh hoạt thay đổi độ phản hồi dựa trên diễn biến cuộc trò chuyện để đáp ứng kì vọng từ khách hàng.

### 2.1.3. Các loại Chatbot

#### **Chatbot dựa trên quy tắc:**

Chatbot dựa trên quy tắc là phiên bản đơn giản nhất của chatbot. Công nghệ này cung cấp cho người dùng các nút hoặc menu để tìm kiếm thông tin cụ thể. Người dùng trải qua một loạt các bước và các câu hỏi được định sẵn để giải quyết vấn đề của họ. Họ không thể nhập vào một câu hỏi mà chỉ có thể nhấp vào một câu hỏi từ một bộ câu hỏi định sẵn. Chatbot có một từ điển tích hợp sẵn để ánh xạ câu trả lời cụ thể cho mọi câu hỏi nên các câu trả lời sẽ giống nhau cho tất cả người dùng dùng về một câu hỏi cụ thể. Chatbot dựa trên quy tắc không phải là lựa chọn tốt cho các tình huống bao gồm các yếu tố chưa xác định, chúng khó điều chỉnh quy mô và mất nhiều thời gian hơn mong muốn để trả lời yêu cầu của người dùng

#### **Chatbot dựa trên từ khóa:**

Các chatbot dựa trên từ khóa hoặc chatbot khai báo sẽ trích xuất các từ khóa cụ thể từ cuộc trò chuyện và cung cấp các câu trả lời tương ứng. Chúng sử dụng các kỹ thuật nhận dạng từ khóa để trích xuất ý định, chủ đề và cảm xúc từ các câu hỏi và trả lời bằng cách sử dụng câu trả lời bằng cách sử dụng câu trả lời theo kịch bản theo những cách được định sẵn. Chatbot loại này vẫn bị hạn chế về phản hồi và chỉ hoạt động trong phạm vi của các chủ đề đã được lập trình sẵn.

#### **Chatbot dựa trên AI:**

Chatbot loại này được áp dụng công nghệ xử lý ngôn ngữ tự nhiên(NLP), hiểu ngôn ngữ tự nhiên(NLU), tạo ngôn ngữ tự nhiên(NLG). AI Tạo sinh cũng giúp chatbot có nhiều khả năng hơn, chatbot được hỗ trợ bởi mô hình ngôn ngữ lớn(LLM) sẽ giúp chatbot mô phỏng lại cuộc trò chuyện tự nhiên.

Các chatbot dựa trên AI tạo sinh cũng có thể xử lý các câu hỏi phức tạp và phát hiện chính xác các cảm xúc cùng thay đổi trong diễn biến trò chuyện và có thể chuyển đổi liên mạch giữa các chủ đề và phản hồi với các cách trả lời phù hợp hơn.

## **2.2. Item-based Recommendation System**

### **2.2.1. Khái niệm**

Hệ thống gợi ý dựa trên mặt hàng (Item-based Recommendation System) là một phương pháp phổ biến trong hệ thống gợi ý, thường được sử dụng để đề xuất các sản phẩm hoặc dịch vụ dựa trên sự tương đồng giữa các mặt hàng.

Nguyên lý hoạt động:

- Hệ thống tập trung vào việc tìm kiếm sự tương đồng giữa các mặt hàng dựa trên hành vi hoặc sở thích của người dùng (Đánh giá, lượt xem, lượt mua,...).
- Nếu 2 mặt hàng có nhiều người dùng tương tác giống nhau (Cùng thích, cùng mua), chúng sẽ được coi là tương đồng.
- Khi người dùng tương tác với một mặt hàng, hệ thống sẽ gợi ý các mặt hàng tương đồng với mặt hàng đó.

Các bước thực hiện của thuật toán:

- Thu thập dữ liệu: Dữ liệu thường là ma trận người dùng – mặt hàng (User – item matrix), trong đó các ô thể hiện tương tác (Đánh giá, lượt xem,...)
- Tính toán độ tương đồng: Sử dụng Cosine Similarity (Góc giữa 2 vector mặt hàng), Pearson Correlation (Đo mối quan hệ tuyến tính giữa các đánh giá của mặt hàng), Jaccard Similarity (Đo mức độ trùng lặp của tập hợp người dùng tương tác).
- Xếp hạng và gợi ý: Chọn các mặt hàng có độ tương đồng cao với mặt hàng mà người dùng đã tương tác và đề xuất chúng

Ưu điểm:

- Hiệu quả với dữ liệu thưa thớt: Dữ liệu nhiều người dùng nhưng ít tương tác.

- Khả năng mở rộng: Dễ tính toán độ tương đồng giữa các mặt hàng trước và lưu trữ.
- Ít bị ảnh hưởng bởi vấn đề người dùng mới (cold-start problem) so với phương pháp dựa trên người dùng.

Nhược điểm:

- Phụ thuộc vào chất lượng dữ liệu: Nếu dữ liệu tương tác ít thì gợi ý có thể không chính xác.
- Thiếu tính cá nhân hóa sâu: Chỉ dựa trên sự tương đồng mặt hàng, không xem xét sâu về sở thích riêng của từng người dùng.
- Khó xử lý các mặt hàng mới: nếu mặt hàng chưa có tương tác thì không thể gợi ý.

So sánh với User-based Recommendation:

Item-based: Tìm mặt hàng tương đồng dựa trên hành vi của tất cả người dùng. Thường nhanh hơn và ổn định hơn.

User-based: Tìm người dùng có sở thích tương tự và gợi ý mặt hàng họ thích, phù hợp khi cần cá nhân hóa nhưng tốn tài nguyên hơn.

### 2.2.2. Cosine Similarity

Cosine Similarity đo lường cosine của góc giữa 2 vector trong không gian đa chiều. Góc càng nhỏ (Gần 0 độ), 2 vector càng tương đồng. Nghĩa là 2 mặt hàng có đặc điểm hoặc hành vi người dùng tương tự.

Trong hệ thống gợi ý, vector thường đại diện cho đánh giá của người dùng đối với một mặt hàng từ ma trận người dùng – mặt hàng.

Giá trị  $\cos \alpha$  nằm trong khoảng  $[-1, 1]$

- Gần 1: 2 mặt hàng hoàn toàn tương đồng.
- Gần -1: 2 mặt hàng có xu hướng đánh giá ngược nhau (đối lập).
- Bằng 0: Không có mối quan hệ tương đồng giữa 2 mặt hàng.

Công thức Cosine Similarity cho 2 mặt hàng

Cho 2 vector A và B đại diện cho 2 mặt hàng, Cosine Similarity giữa 2 mặt hàng được tính như sau:

$$\text{Cosine Similarity}(A, B) = \frac{\sum_{i=1}^n A_i \cdot B_i}{\sqrt{\sum_{i=1}^n A_i^2} \cdot \sqrt{\sum_{i=1}^n B_i^2}}$$

Công thức Cosine Similarity trong hệ thống gợi ý:

Trong hệ thống gợi ý,  $A_i$  và  $B_i$  thường là đánh giá của người dùng  $i$  cho 2 mặt hàng  $A$  và  $B$ . Tuy nhiên không phải mọi người dùng đều đánh giá cả 2 mặt hàng nên ta chỉ tính tổng trên tập hợp người dùng đã đánh giá cả  $A$  và  $B$ .

$$\text{Cosine Similarity}(A, B) = \frac{\sum_{u \in U_{A \cap B}} r_{u,A} \cdot r_{u,B}}{\sqrt{\sum_{u \in U_{A \cap B}} r_{u,A}^2} \cdot \sqrt{\sum_{u \in U_{A \cap B}} r_{u,B}^2}}$$

Trong đó:

- $U_{A \cap B}$ : Tập hợp người dùng đã đánh giá cả 2 mặt hàng  $A$  và  $B$ .
- $r_{u,A}$ : Đánh giá của người dùng  $u$  cho mặt hàng  $A$ .
- $r_{u,B}$ : Đánh giá của người dùng  $u$  cho mặt hàng  $B$ .

#### Adjust Cosine Similarity

Trong hệ thống gợi ý, Cosine Similarity cơ bản có thể bị ảnh hưởng bởi sự khác biệt trong cách người dùng đánh giá (Một người cho điểm cao trong khi người kia cho điểm thấp). Adjust Cosine Similarity sẽ giải quyết vấn đề bằng cách chuẩn hóa đánh giá của người dùng bằng cách trừ đi điểm đánh giá trung bình của họ.

Công thức:

$$\text{Adjusted Cosine Similarity}(A, B) = \frac{\sum_{u \in U_{A \cap B}} (r_{u,A} - \bar{r}_u) \cdot (r_{u,B} - \bar{r}_u)}{\sqrt{\sum_{u \in U_{A \cap B}} (r_{u,A} - \bar{r}_u)^2} \cdot \sqrt{\sum_{u \in U_{A \cap B}} (r_{u,B} - \bar{r}_u)^2}}$$

Với  $\bar{r}_u$  là điểm đánh giá trung bình của người dùng  $u$  trên tất cả các mặt hàng họ đã đánh giá.

## 2.3. Logistic Regression (LR)

### 2.3.1. Khái niệm

Hồi quy Logistic là thuật toán học có giám sát để cho những nhiệm vụ phân loại với mục đích dự đoán xác suất của một cá thể thuộc về một lớp nhất định hay không. Hồi quy logistic là một thuật toán thống kê phân tích mối quan hệ giữa một hoặc nhiều biến độc lập với một biến phụ thuộc phân loại.

Hồi quy Logistic được sử dụng để phân loại nhị phân, sử dụng hàm sigmoid. Lấy đầu vào làm biến độc lập và tạo ra giá trị xác suất từ 0 tới 1.

VD: Có hai lớp 0 và lớp 1, nếu xác suất dự đoán từ hàm Sigmoid lớn hơn 0.5, mẫu đó được phân vào lớp 1, ngược lại là lớp 0.

Đặc điểm chính:

- Hồi quy logistic dự đoán đầu ra của một biến phụ thuộc phân loại. Kết quả phải là một giá trị phân loại hoặc rời rạc.
- Nó có thể là Có hoặc 0, 0 hoặc 1, đúng hoặc sai. Tuy nhiên thay vì đưa ra giá trị chính xác là 0 và 1, nó đưa ra các giá trị xác suất nằm trong khoảng từ 0 đến 1.
- Trong hồi quy Logistic, ta sử dụng hàm Sigmoid để dự đoán giá trị.

Các loại hồi quy logistic:

- Nhị thức: Chỉ có thể có 2 loại biến phụ thuộc, VD: 0 hoặc 1, Đạt hoặc không đạt
- Đa thức: Có thể có 3 hoặc nhiều loại biến phụ thuộc không theo thứ tự. VD: Chó, mèo, cừu,...
- Thứ tự: Có thể có 3 hoặc nhiều loại biến phụ thuộc có trật tự. VD: Thấp, Trung Bình, Cao.

Giả định của hồi quy logistic:

- Quan sát độc lập: Mỗi quan sát độc lập với quan sát còn lại. Không có mối tương quan giữa bất kỳ biến đầu vào nào.
- Biến phụ thuộc nhị phân: Biến phụ thuộc phải là nhị phân, chỉ có thể lấy 2 giá trị
- Mối quan hệ giữa các biến độc lập và các log của các biến phụ thuộc phải là tuyến tính.
- Kích thước mẫu đủ lớn.

Sigmoid:

- Hàm sigmoid là một hàm toán học sử dụng để ánh xạ các giá trị dự đoán với xác suất.
- Ánh xạ bất kỳ giá trị thực nào vào một giá trị khác trong phạm vi 0 và 1. Giá trị của hồi quy logistic phải nằm trong khoảng 0 đến 1. Nên nó sẽ tạo thành một đường cong giống như dạng “S”.
- Đường cong dạng S gọi là hàm Sigmoid hoặc hàm logistic.



- Giá trị ngưỡng phổ biến nhất là 0.5, tức là nếu xác suất lớn hơn 0.5 thì thuộc lớp 1, ngược lại thuộc lớp 0.

Cách hoạt động của hồi quy logistic:

Mô hình hồi quy logistic chuyển đổi đầu ra giá trị liên tục của hàm hồi quy tuyến tính thành đầu ra giá trị phân loại bằng cách sử dụng hàm sigmoid, ánh xạ bất kỳ tập hợp đầu vào các biến độc lập có giá trị thực nào thành giá trị từ 0 đến 1.

Thử cho các tính năng đầu vào độc lập là:

$$X = \begin{matrix} x_{11} & \cdots & x_{1m} \\ x_{21} & \cdots & x_{2m} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{Nm} \end{matrix}$$

Biến phụ thuộc là Y chỉ có giá trị nhị phân, tức là 0 hoặc 1.

$$Y = \begin{cases} 0 & \text{nếu Class1} \\ 1 & \text{nếu Class2} \end{cases}$$

Sau đó áp dụng hàm đa tuyến tính cho các biến đầu vào X:

$$z = \left( \sum_{i=1}^n w_i x_i \right) + b$$

Ở đây  $x_i$  là quan sát thứ i,  $w_i = [w_1, w_2, w_3, \dots, w_m]$  là trọng số hoặc hệ số hồi quy (coefficients).  $b$  là hệ số chặn (bias) hay hệ số tự do (intercept). Công thức tổng quát của mô hình hồi quy tuyến tính có thể biểu diễn dưới dạng tích vô hướng.

Công thức tổng quát của mô hình hồi quy tuyến tính có thể biểu diễn dưới dạng tích vô hướng như sau:

$$z = w \cdot X + b$$

$z$ : Giá trị dự đoán (đầu ra dự đoán của mô hình hồi quy tuyến tính).

$w$ : Trọng số (Quyết định mức độ ảnh hưởng của từng đặc trưng X đến giá trị dự đoán  $z$ ).

$X$ : Tập dữ liệu đầu vào, gồm các biến độc lập (features) mà mô hình sử dụng để dự đoán giá trị đầu ra. Nếu có nhiều biến đầu vào,  $X$  là một vector chứa nhiều giá trị.

$b$ : Hệ số chặn (Giá trị điều chỉnh giúp mô hình khớp tốt hơn với dữ liệu thực tế, xác định điểm mà đường hồi quy cắt trục tung khi tất cả biến đầu vào  $X$  bằng 0).

Hàm Sigmoid được định nghĩa như sau:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

$z$ : Giá trị đầu vào từ hồi quy tuyến tính.

$e$ : Hằng số Euler ( $\sim 2.718$ )

Hàm sigmoid giúp chuyển đổi giá trị  $z$  (có thể nằm trong khoảng  $(-\infty, +\infty)$ ) thành giá trị xác suất từ 0 đến 1.

Tính chất của hàm sigmoid:

- Khi  $z \rightarrow +\infty, \sigma(z) \rightarrow 1$  : Xác suất rất cao cho lớp 1.
- Khi  $z \rightarrow -\infty, \sigma(z) \rightarrow 0$  : Xác suất rất thấp cho lớp 1.
- Khi  $z = 0, \sigma(0) = 0.5$ : Xác suất trung bình (Không thiên về lớp nào)

Ta sử dụng Sigmoid để tính xác suất một điểm dữ liệu thuộc về lớp 1 ( $y = 1$ ) như sau:

$$P(y = 1) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

$$P(y = 0) = 1 - \sigma(z)$$

$P(y = 1)$ : Là xác suất để mẫu thuộc lớp 1.

$P(y = 0)$ : Là xác suất để mẫu thuộc lớp 0.

Mô hình hồi quy Logistic sử dụng hàm Sigmoid để tính xác suất một điểm dữ liệu thuộc về một lớp cụ thể, từ đó phân loại dữ liệu dựa trên ngưỡng (thường là 0.5).

Biến đổi thành phương trình hồi quy logistic:

Hồi quy logistic sử dụng khái niệm odds ratio: Tức là tỷ lệ một sự kiện xảy ra so với không xảy ra:

$$Odds = \frac{p(x)}{1 - p(x)}$$

Trong mô hình hồi quy Logistic, ta giả định rằng odds ratio có dạng hàm mũ:

$$\frac{p(x)}{1 - p(x)} = e^z$$

Với  $z = wX + b$

Lấy logarit tự nhiên 2 vế, ta có:

$$\ln\left(\frac{p(x)}{1 - p(x)}\right) = z = w \cdot X + b$$

Đây là phương trình logit, biểu diễn quan hệ tuyến tính giữa log-odds và dữ liệu đầu vào. Bằng cách triển khai phương trình, ta thu được công thức cuối cùng của hồi quy logistic:

$$p(X; w, b) = \frac{e^{w \cdot X + b}}{1 + e^{w \cdot X + b}}$$

Hay:

$$p(X; w; b) = \frac{1}{1 + e^{-(wX+b)}}$$

Phương trình logit cho thấy log-odds phụ thuộc tuyến tính vào dữ liệu đầu vào. Hàm Sigmoid giúp chuyển đổi giá trị log-odds sang xác suất trong khoảng (0, 1).

### Hàm khả năng (Likelihood Function) trong hồi quy Logistic

Trong hồi quy logistic, xác suất dự đoán cho mỗi quan sát được tính bằng:

$$\text{Khi } y_i = 1: p(x_i; b, w) = p_i$$

$$\text{Khi } y_i = 0: 1 - p(x_i; b, w) = 1 - p_i$$

Tức là:

$$P(y_i | x_i; w, b) = p_i^{y_i} (1 - p_i)^{(1-y_i)}$$

Đây là xác suất có điều kiện của nhãn  $y_i$  cho một điểm dữ liệu  $x_i$

Với hàm sigmoid:

$$p_i = \frac{1}{1 + e^{-(wx_i+b)}}$$

Hàm khả năng (Likelihood Function) biểu diễn xác suất của toàn bộ tập dữ liệu (gồm n quan sát) dựa trên mô hình hồi quy logistic:

$$L(b, w) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{(1-y_i)}$$

Với:

$$p_i = \frac{1}{1 + e^{-(wx_i+b)}}$$

$L(b, w)$ : Thể hiện mức độ phù hợp của mô hình với dữ liệu quan sát.

Hàm log-khả năng (Log-likelihood)

Lấy log 2 về để chuyển tích thành tổng:

$$\text{Log}L(b, w) = \sum_{i=1}^n y_i \log p(x_i) + (1 - y_i) \log(1 - p(x_i))$$

Thay thế  $p(x_i) = \frac{1}{1 + e^{-(w.x_i + b)}}$  ta có:

$$\text{Log}L(b, w) = \sum_{i=1}^n y_i (wx_i + b) - \sum_{i=1}^n \log(1 + e^{wx_i + b})$$

Hàm log-khả năng trên là mục tiêu tối đa hóa trong hồi quy logistic.

Gradient của hàm log-khả năng:

Ta cần tính đạo hàm theo  $w$  để tìm cực đại của  $\text{LogL}(b, w)$ :

$$\frac{\partial \log L(b, w)}{\partial w_j} = \sum_{i=1}^n (y_i - p(x_i; b, w)) x_{ij}$$

Tương tự đạo hàm theo  $b$ :

$$\frac{\partial \text{LogL}(b, w)}{\partial b} = \sum_{i=1}^n (y_i - p(x_i; b, w))$$

Biểu thức thể hiện rằng chênh lệch giữa giá trị thực  $y_i$  và xác suất dự đoán  $p(x_i)$  sẽ quyết định sự thay đổi của trọng số  $w$ .

Khi giá trị dự đoán  $p(x_i)$  sai lệch nhiều so với  $y_i$ , ta cần cập nhật  $w$  để giảm sai số

Ta cập nhật trọng số theo công thức ( $\alpha$  là tốc độ học) :

$$w_j = w_j + \alpha \sum_{i=1}^n (y_i - p(x_i)) x_{ij}$$

### 2.3.2. Multinomial Logistic Regression

Trong phân loại đa lớp, Logistic Regression sử dụng hàm Softmax thay vì Sigmoid để tính xác suất thuộc vào từng lớp.

Cho tập dữ liệu:

- $x = [x_1, x_2, \dots, x_n]$ : Vector đặc trưng của một mẫu (Ví dụ: Các từ đã được mã hóa).
- $y$ : Nhãn lớp (Ví dụ:  $y = 0$  cho tiêu cực,  $y = 1$  cho trung tính,  $y = 2$  cho tích cực).
- $K = 3$ : Số lớp (Tiêu cực, trung tính, tích cực).
- $w_k$ :  $[w_{k1}, w_{k2}, \dots, w_{kn}]$ : Vector trọng số cho lớp  $k$ .
- $b_k$ : Hằng số bias cho lớp  $k$

Mô hình tính điểm số cho mỗi lớp  $k$  bằng tổ hợp tuyến tính:

$$z_k = w_k \cdot x + b_k = w_{k1}x_1 + \dots + w_{kn}x_n + b_k$$

Xác suất mẫu thuộc lớp  $k$  được tính bằng hàm Softmax:

$$P(y = k|x) = \frac{e^{z_k}}{\sum_{j=0}^{K-1} e^{z_j}}$$

$e^{z_k}$ : Chuyển đổi điểm số của lớp  $k$  thành dạng exponent để đảm bảo giá trị dương.

$\sum_{j=0}^{K-1} e^{z_j}$ : Tổng các giá trị exponent của tất cả các lớp dùng để chuẩn hóa xác suất (Tổng xác suất của các lớp bằng 1)

Để phân loại, mô hình chọn lớp có xác suất cao nhất.

Để huấn luyện mô hình, Multinomial Logistic Regression sử dụng hàm mất mát Cross-Entropy Loss (Log loss) cho phân loại đa lớp.

Hàm mất mát cho một mẫu

Cho một mẫu  $(x, y)$  trong đó  $y = k$  là lớp thực tế, hàm mất mát được tính:

$$Loss = -\log \left( \frac{e^{z_k}}{\sum_{j=0}^{K-1} e^{z_j}} \right)$$

Với  $N$  mẫu, thì hàm mất mát trung bình là:

$$J(w, b) = -\frac{1}{N} \sum_{i=1}^N \sum_{k=0}^{K-1} 1\{y_i = k\} \log (P(y_i = k|x_i))$$

Trong đó:

- $1\{y_i = k\}$ : Bằng 1 nếu mẫu  $i$  thuộc lớp  $k$ , bằng 0 nếu không.
- $P(y_i = k|x_i)$ : Xác suất dự đoán cho lớp  $k$  của mẫu  $i$ .

Mục tiêu là tối ưu hóa  $J(w, b)$  bằng cách điều chỉnh  $w_k$  và  $b_k$  thông qua Gradient Descent hoặc các phương pháp tối ưu khác.

## PHẦN 3: PHÂN TÍCH, THIẾT KẾ

### 3.1. Ý tưởng thiết kế

Trọng tâm của thiết kế là sự phối hợp giữa các phương pháp phân tích ngôn ngữ tự nhiên và học máy để nhận biết trạng thái tâm lý từ văn bản đầu vào. Các kỹ thuật như phân tích từ vựng dựa trên luật, tương tự cách tiếp cận của VADER, được sử dụng để đánh giá mức độ tích cực, tiêu cực hoặc trung lập của nội dung giao tiếp, chú trọng đến các yếu tố như từ ngữ, ngữ cảnh, và dấu câu. Đồng thời, các mô hình học máy, chẳng hạn như hồi quy Logistic đa lớp, được áp dụng để phân loại cảm xúc một cách chính xác hơn, tận dụng dữ liệu huấn luyện để nhận diện các mẫu cảm xúc phức tạp. Những phương pháp này được kết hợp để đảm bảo rằng hệ thống không chỉ nhận biết trạng thái tâm lý mà còn hiểu được sự chuyển đổi tinh tế trong cảm xúc, từ đó tạo ra các phản hồi phù hợp, ví dụ như an ủi khi người dùng tỏ ra thất vọng hoặc khuyến khích khi họ thể hiện sự hứng thú.

Song song với phân tích cảm xúc, hệ thống tích hợp lý thuyết về hệ khuyến nghị dựa trên sự tương đồng giữa các mặt hàng, sử dụng các kỹ thuật như Cosine Similarity để đo lường mức độ tương đồng dựa trên dữ liệu hành vi người dùng. Bằng cách phân tích các mẫu tương tác, hệ thống xác định các sản phẩm hoặc dịch vụ có liên quan chặt chẽ, từ đó đưa ra gợi ý phù hợp với sở thích cá nhân. Sự kết hợp giữa phân tích cảm xúc và khuyến nghị được thiết kế để hoạt động đồng bộ: trạng thái tâm lý của người dùng sẽ ảnh hưởng đến cách các gợi ý được trình bày, chẳng hạn như ưu tiên các sản phẩm mang tính tích cực khi người dùng vui vẻ. Các lý thuyết này được tích hợp trong một kiến trúc linh hoạt, cho phép hệ thống xử lý dữ liệu thời gian thực và tạo ra phản hồi mang tính cá nhân hóa cao. Thiết kế nhấn mạnh tính mô-đun, đảm bảo rằng mỗi thành phần lý thuyết có thể được tinh chỉnh độc lập mà vẫn duy trì sự thống nhất, đồng thời cho phép mở rộng sang các ứng dụng mới như tư vấn học tập hoặc hỗ trợ sức khỏe tâm lý, mang lại giá trị lâu dài trong việc nâng cao chất lượng tương tác người-máy.

## 3.2. Công nghệ được áp dụng

### 3.2.1. Django

Django là một framework web được viết bằng ngôn ngữ lập trình Python được thiết kế để xây dựng ứng dụng web mạnh mẽ, bảo mật và dễ mở rộng.

Django sử dụng mô hình MTV (Model-Template-View), một biến thể của MVC (Model-View-Controller):

- Model: Quản lý dữ liệu và logic nghiệp vụ (Tương tác với cơ sở dữ liệu).
- Template: Xử lý giao diện người dùng (HTML, CSS, JavaScript).
- View: Xử lý logic hiển thị và tương tác với người dùng.

Các đặc điểm chính của Django:

- ORM (Object-Relational Mapping): Django cung cấp một ORM mạnh mẽ, cho phép người dùng tương tác với cơ sở dữ liệu (Như PostgreSQL, MySQL, SQLite) bằng cú pháp Python thay vì SQL trực tiếp.
- Hệ thống quản trị (Admin Interface): Django tự động tạo giao diện quản trị (admin panel) cho phép quản lý dữ liệu (thêm, sửa, xóa) mà không cần viết mã bổ sung.
- Hệ thống URL Routing: Django sử dụng hệ thống định tuyến URL dựa trên biểu thức chính quy, giúp ánh xạ các yêu cầu HTTP đến các hàm xử lý (views) một cách linh hoạt.
- Template Engine: Cho phép tách biệt logic nghiệp vụ và giao diện người dùng, template hỗ trợ các thẻ và bộ lọc để hiển thị dữ liệu động.
- Bảo mật tích hợp: Django cung cấp các tính năng bảo mật: Ngăn chặn tấn công Cross-Site Request Forgery, bảo vệ tấn công Cross-Site Scripting, Phòng tránh SQL Injection, Authentication và Authorization.
- Hỗ trợ đa ngôn ngữ và quốc tế hóa: Hỗ trợ dịch ứng dụng sang nhiều ngôn ngữ và xử lý múi giờ, phù hợp với các ứng dụng toàn cầu.

Cấu trúc một dự án Django:

myproject/ ├── manage.py v.v.) └── myproject/	# Tập điều khiển các lệnh Django (runserver, migrate,
--	---

—	__init__.py	
—	settings.py	# Cấu hình dự án (database, app, middleware, v.v.)
—	urls.py	# Định tuyến URL
—	wsgi.py	# Tập cấu hình triển khai WSGI
	myapp/	
—	__init__.py	
—	admin.py	# Cấu hình giao diện admin
—	apps.py	# Cấu hình ứng dụng
—	migrations/	# Lưu trữ các tệp di cư cơ sở dữ liệu
—	models.py	# Định nghĩa mô hình dữ liệu
—	tests.py	# Viết unit test
—	views.py	# Xử lý logic hiển thị
—	templates/	# Lưu trữ tệp template (HTML)
—	static/	# Lưu trữ tệp tĩnh (CSS, JS, hình ảnh)

### 3.2.2. VADER

VADER (Valence Aware Dictionary and Sentiment Reasoner) là một công cụ phân tích cảm xúc dựa trên từ điển và luật (lexicon-based sentiment analysis) được thiết kế đặc biệt để xử lý các văn bản ngắn, chẳng hạn như bài đăng trên mạng xã hội, bình luận hoặc đánh giá. VADER được sử dụng rộng rãi trong các bài toán phân loại cảm xúc (sentiment classification) vì tính đơn giản, hiệu quả và khả năng xử lý ngôn ngữ tự nhiên trong các ngữ cảnh không chính thức. Mục đích của VADER là phân tích cảm xúc của văn bản, xác định xem văn bản thể hiện cảm xúc tích cực, tiêu cực hay trung tính, đồng thời cung cấp mức độ cảm xúc (intensity) thông qua điểm số.

Cách hoạt động của VADER:

VADER sử dụng một từ điển chứa khoảng 7500 từ, cụm từ, biểu tượng cảm xúc và các thành phần ngôn ngữ, mỗi thành phần được gán một điểm valance (Độ tích cực hoặc tiêu cực).

Điểm valance nằm trong khoảng [-4,4]:

- +4: Cực kỳ tích cực.
- -4: Cực kỳ tiêu cực.
- 0: Trung tính hoặc không mang cảm xúc.

VADER không chỉ dựa vào điểm từ điển mà còn áp dụng các luật ngữ pháp và ngữ cảnh để điều chỉnh điểm số cảm xúc. Các luật này bao gồm:



- **Tăng cường hoặc giảm cường độ:** Từ ngữ tăng cường (intensifiers) như “very”, “extremely” làm tăng điểm cảm xúc (Ví dụ “verygood” sẽ có điểm là +2.0, tích cực hơn “good” với số điểm là +1.3). Ngược lại từ ngữ giảm nhẹ như “slightly”, “kind of” làm giảm điểm (Ví dụ “slightly bad” sẽ có điểm là -0.8 ít tiêu cực hơn “bad” với -1.5 điểm).
- **Viết hoa và dấu câu:** Viết hoa làm tăng cường độ cảm xúc (Ví dụ như “GREAT” có +2.5 điểm sẽ tích cực hơn “great” với +1.8 điểm). Dấu chấm than(!) hoặc dấu (?) cũng sẽ tăng cường độ (Ví dụ như “This is bad!” Có -2 điểm tiêu cực hơn “This is bad” có -1.5 điểm).
- **Phủ định:** Các từ phủ định như “not”, “never”, “no” đảo ngược hoặc giảm cảm xúc.
- **Biểu tượng cảm xúc và từ lóng:** VADER nhận diện các emojis, emoticons và các từ lóng phổ biến.
- **Từ nối:** Ví dụ như từ but giảm các xúc của phần trước và tăng cảm xúc của phần sau (Ví dụ: “The movie was good but boring sẽ ưu tiên cảm xúc của boring).

Đầu ra của VADER là một điểm số cảm xúc cho văn bản, bao gồm:

- **Điểm riêng lẻ cho từng loại cảm xúc:**
  - Tích cực: Tỷ lệ cảm xúc tích cực.
  - Tiêu cực: Tỷ lệ cảm xúc tiêu cực.
  - Trung tính: Tỷ lệ cảm xúc trung tính.
  - Tổng của Tích cực + Tiêu cực + Trung tính = 1.0.
- **Điểm tổng hợp (Compound score):** Là một giá trị tổng hợp trong khoảng [-1;1]:
  - $> 0$ : Tích cực
  - $< 0$ : Tiêu cực
  - $\approx 0$ : Trung tính

Điểm compound được chuẩn hóa từ tổng các điểm valance của các từ trong câu, có tính đến các luật Heuristic.

Dựa trên điểm Compound ta áp dụng các ngưỡng để phân loại:

- Tích cực:  $compound \geq 0.05$
- Trung tính:  $-0.05 < compound < 0.05$ .

- Tiêu cực:  $compound \leq -0.05$ .

### 3.3. Tài nguyên được sử dụng

Link Source Code: <https://github.com/HueyAnthonyDisward/Chatbot-with-recomendation-system-and-sentiment-analyst>

#### 3.3.1. Tập dữ liệu

##### 3.3.1.1. Bitext Gen AI Chatbot Customer Support Dataset

Link Dataset: <https://www.kaggle.com/datasets/bitext/bitext-gen-ai-chatbot-customer-support-dataset>

Mô tả Dataset: Dataset này dùng để huấn luyện các mô hình ngôn ngữ lớn (LLM) như GPT, Llama2, Falcon cho Fine Tuning và Domain Adaption trong lĩnh vực dịch vụ khách hàng. Dataset này dùng để huấn luyện mô hình để phát hiện ý định (Intent Detection). Với quy mô gồm 26.872 cặp câu hỏi/trả lời, khoảng 1000 cặp cho mỗi ý định.

Các danh mục và các intent:

- Danh mục: Account, cancellation\_fee, contact, delivery, feedback, invoice, order, payment, refund, subscription.
- Intent: Create\_account, check\_refund\_policy,...

*Bảng 3.3.1.1.1: Danh Sách Cột Và Thuộc Tính*

Tên cột	Tên thuộc tính
flags	Thẻ biến thể ngôn ngữ
instruction	Yêu cầu của người dùng
category	Danh mục các intent
intent	Ý định cụ thể
response	Phản hồi mẫu

##### 3.3.1.2. Consumer Reviews of Amazon Products

Link Dataset: <https://www.kaggle.com/datasets/datafiniti/consumer-reviews-of-amazon-products>

Dataset chứa hơn 34.000 đánh giá của người tiêu dùng về các sản phẩm của Amazon như Kindle, Fire TV và các sản phẩm điện tử khác. Dataset này phù hợp với các ứng dụng phân tích cảm xúc, hoặc nghiên cứu hành vi người tiêu dùng.

Các thuộc tính trong dataset:

*Bảng 3.3.1.2.1: Danh Sách Cột Và Ý Nghĩa*

Tên cột	Ý nghĩa
id	Mã định danh duy nhất của sản phẩm.
name	Tên của sản phẩm
asins	Mã định danh của Amazon
brand	Thương hiệu sản phẩm
categories	Danh mục
keys	Từ khóa liên quan đến các sản phẩm
manufacturer	Nhà sản xuất của sản phẩm
reviews.date	Ngày người dùng viết đánh giá
reviews.dateAdded	Ngày đánh giá được thêm vào hệ thống
reviews.dateSeen	Ngày hệ thống cuối cùng thấy đánh giá này
reviews.didPurchase	Người dùng có thực sự mua sản phẩm hay không.
reviews.doRecommend	Người dùng có đề xuất sản phẩm này không
reviews.id	Mã định danh của đánh giá
reviews.numHelpful	Số người thấy đánh giá hữu ích
reviews.rating	Điểm đánh giá
reviews.sourceURLs	URL nguồn nơi lấy đánh giá
reviews.text	Nội dung đánh giá của người dùng
reviews.title	Tiêu đề của đánh giá
reviews.userCity	Thành phố của người đánh giá
reviews.username	Tên người dùng

### 3.3.2. Thư mục

#### 3.3.2.1. Thư mục Dataset

Đây là thư mục chứa các dataset và để xử lý data trước khi đưa vào mô hình.

##### **Xử lý dataset cho hệ khuyến nghị:**

Dataset gốc về đánh giá các sản phẩm Amazon chứa rất nhiều cột thuộc tính. Nên để đưa vào mô hình, chúng ta cần tiền xử lý dữ liệu trước: Các bước tiền xử lý gồm:

- Giữ lại các cột liên quan: Asins (Mã sản phẩm), name (Tên sản phẩm), categories (Danh mục sản phẩm), reviews.text (Nội dung đánh giá), reviews.rating (Điểm đánh giá).
- Xóa các hàng có giá trị Nan trong bất kỳ cột nào trong các cột liên quan, bỏ đi các hàng trùng lặp.
- Chuẩn hóa chuỗi: Loại bỏ dấu ngoặc vuông ở đầu và cuối chuỗi, tách chuỗi thành danh sách dựa trên dấu phẩy và lấy các danh mục.
- In số lượng sản phẩm thuộc từng danh mục.
- Chọn 10 danh mục có số lượng sản phẩm lớn nhất và đưa vào List.
- Lọc các sản phẩm dựa trên các danh mục.
- Lưu dữ liệu đã làm sạch.

Kết quả, ta sẽ được một file dataset mới có tên là `cleaned_amazon_reviews.csv`.

##### **Xử lý dataset cho chatbot:**

Để huấn luyện cho chatbot, ta cần tập trung vào 2 cột thuộc tính trong dataset, đó là thuộc tính intent (Ý định) và thuộc tính Response (phản hồi), quy trình tiền xử lý dataset Bibtex như sau:

- Kiểm tra giá trị duy nhất trong category và intent để có thể hiểu được nội dung của dataset.
- Loại bỏ các dòng bị thiếu.
- Xóa cột flags không cần thiết.
- Chuẩn hóa dữ liệu bằng cách chữ hoa được chuyển thành chữ thường trong intent để tránh trùng lặp, ép kiểu tất cả về dạng str để đảm bảo đồng nhất.
- Lưu lại tập dữ liệu đã được làm sạch.

##### **Huấn luyện cho Chatbot:**

Chatbot trong project được huấn luyện theo mô hình phân loại intent (Ý định của người dùng) từ tập dữ liệu bitext đã được xử lý qua. Sử dụng Logistic Regression và TfidfVectorizer để huấn luyện.

Mô hình sẽ được huấn luyện với tập các dữ liệu đầu vào là các câu hỏi hoặc câu lệnh từ người dùng và nhãn đầu ra chính là intent tức là ý định của người dùng.

Mô hình sử dụng TF-IDF để chuyển câu văn bản thành vector số để đưa vào mô hình học máy, sử dụng thuật toán LogisticRegression để phân loại.

Cuối cùng, đầu ra của chương trình sẽ là file intent\_classifier.pkl để đưa vào sử dụng cho chatbot.

### **Placeholder:**

Placeholder là giá trị tạm thời dùng để hiển thị hoặc đánh dấu vị trí mà nội dung thật sẽ được điền vào sau. Trong tập dữ liệu bitext có sự xuất hiện rất nhiều của các Placeholder nhằm để có thể tùy chỉnh cho đa dạng chủ đề. Ta cần làm một file để hiển thị Placeholder có trong dataset để hiểu về dữ liệu cũng như xác định được chiến lược thiết kế trong tương lai. Quy trình xác định placeholder như sau:

- Xác định các placeholder: Các Placeholder trong dataset sẽ có dạng 2 cặp ngoặc kép{{<Nội dung>}}.
- Tìm các placeholder duy nhất, tạo một set để đưa các placeholders vào mà không trùng lặp.
- Sắp xếp lại và in ra kết quả.

### **3.3.2.2. Thư mục Sentiment**

Đây là thư mục dùng để xử lý dữ liệu và tìm ra các sản phẩm đạt chất lượng cao. Kết hợp giữa TF-IDF cùng với Cosine Similarity với quy trình:

- Phân tích cảm xúc của review: Sử dụng VADER để tính điểm cảm xúc với ngưỡng compound  $\geq 0.05$  là tích cực, ngược lại là tiêu cực. Phân tích toàn bộ mẫu trong cột reviews.text.
- Sau khi có được tất cả cảm xúc của reviews, ta gộp dữ liệu lại theo asins và tính tỷ lệ phần trăm đánh giá tích cực cho từng sản phẩm. Giữ lại các sản phẩm có tỷ lệ tích cực  $> 80\%$ .

- Xây dựng ma trận TF-IDF cho từng ASIN, gom các review thành một đoạn văn duy nhất cho mỗi ASIN. Chuyển các đoạn văn thành vector TF-IDF và tính ma trận độ tương đồng cosine giữa các sản phẩm. để ra được hàm gợi ý sản phẩm.
- Hàm Gợi ý sản phẩm sẽ lấy chỉ số sản phẩm và tính điểm tương đồng và sắp xếp các sản phẩm tương tự theo điểm cosine giảm dần, lấy các thông tin trong sản phẩm và lọc theo từ khóa.

### 3.3.2.3. Thư mục Collaborative Recommendation

Cùng mục đích với thư mục sentiment, cũng là để xây dựng hệ khuyến nghị nhưng ở Collaborative Recommendation sẽ dựa vào điểm đánh giá chứ không phải dựa vào bình luận của người dùng. Quy trình như sau:

- Ta tạo hồ sơ sản phẩm bằng cách gộp dữ liệu theo các ASIN tức là mỗi ASIN chỉ có một dòng, lấy name và category đầu tiên cho mỗi ASIN cùng với tính điểm đánh giá trung bình cho mỗi sản phẩm.
- Tiếp theo, ta tính điểm đánh giá trung bình cho mỗi sản phẩm gồm ASIN và điểm đánh giá trung bình.
- Sử dụng độ tương đồng cosine giữa các vector average\_rating của sản phẩm.

So sánh 2 phương án xây dựng hệ khuyến nghị:

**Bảng 3.3.2.3.1: Hai Phương Án Xây Dựng Hệ Khuyến Nghị**

Tiêu chí	TF-IDF + Sentiment Analysis	Cosine Similarity dựa trên Rating
Tiền xử lý dữ liệu	Dữ liệu ASIN được giữ nguyên	Tách từng ASIN riêng biệt nếu có nhiều ASIN trong một dòng
Sentiment Analysis	Sử dụng Vader để xác định cảm xúc đánh giá	Không sử dụng
Chọn sản phẩm chất lượng cao	Lọc ra các sản phẩm có tỷ lệ đánh giá tích cực >80%	Không lọc theo đánh giá, tất cả sản phẩm đều được xem xét
Phân tích nội dung	Dùng TF-IDF để biến nội dung review thành vector(5000 từ)	Không dùng nội dung review, chỉ dùng rating trung bình
Cách tính	Tính toán cosine similarity giữa các vector tf-idf (Theo nội dung review)	Tính cosine similarity giữa các vector rating trung bình)

Tính chính xác của gợi ý	Cao hơn vì xét cảm xúc, lọc review chất lượng và dùng nội dung bình luận	Thấp hơn vì chỉ xét rating trung bình(Không đủ chiều thông tin)
Mức độ phù hợp	Phù hợp cho recommendation thực tế dựa trên nội dung đánh giá	Chỉ phù hợp cho gợi ý đơn giản, không xét nội dung review

#### 3.3.2.4. Thư mục Sentiment\_User

Đây là thư mục dùng để nhận diện cảm xúc của người dùng, cũng sử dụng TF-IDF cùng với VADER và Logistic Regression để phân loại cảm xúc của người dùng theo quy trình:

- Gán nhãn cảm xúc với VADER, ta gán nhãn cho từng người, positive nếu  $\text{compound\_score} \geq 0.05$ , negative với  $\text{compound\_score} \leq -0.05$  và neutral nếu nằm giữa. Sau đó ta áp dụng nhãn vào dữ liệu.
- Ta xác định rằng, tập dữ liệu chính là các câu hỏi của người dùng và tập nhãn là các cảm xúc(positive, negative cùng với  $\text{compound\_score}$ ). Ta chia tỷ lệ 80% là tập huấn luyện và 20% là tập kiểm tra.
- Sau đó ta vector hóa văn bản bằng TF-IDF, loại bỏ các stop word và đưa vào mô hình Logistic Regression. Cuối cùng là đánh giá mô hình cùng với lưu mô hình để sử dụng trong Django.

#### 3.3.2.5. Thư mục Web

Đây là bước cuối cùng, sau khi chúng ta đã xử lý dữ liệu, huấn luyện mô hình cùng với việc lưu lại các mô hình để sử dụng thì chúng ta sẽ tích hợp nó với web Django.

Webapp Amazon\_Chatbot\_Project gồm các Logic xử lý tích hợp (Chatbot, hệ khuyến nghị, nhận diện cảm xúc) và giao diện chatbot trên web.

Với thư mục data là thư mục chứa các dataset, mô hình cùng với kết quả vector hóa từ các quy trình trước.

Ở webapp ta cần chú ý đến 2 file đó chính là chatbot.html và views.py:

- Chatbot.html: Đây là file dựng giao diện trò chuyện người dùng với chatbot, hiển thị ảnh đại diện, hiển thị nội dung tin nhắn, phân tích cảm xúc của người dùng và gửi/nhận dữ liệu từ server(endpoint /chat/).

- Views.py: Phần logic cho Chatbot, với chức năng: Dự đoán intent từ mô hình Logistic Regression được huấn luyện từ Bitext, phân tích cảm xúc của nội dung người dùng nhập, đề xuất sản phẩm từ dữ liệu Amazon(Dựa trên TF-IDF và Cosine similarity) cùng với trả lời dựa trên dữ liệu từ bitext.



## PHẦN 4: THỰC NGHIỆM

### 4.1. Kết quả thực nghiệm

Đầu tiên là kết quả của thư mục dataset, thư mục dataset có mục đích lưu trữ, làm sạch dữ liệu, tiền xử lý dữ liệu cùng với chuẩn bị model phân loại ý định:

Name	Date modified	Type	Size
customer-support-llm-chatbot-training-...	9/23/2023 1:55 AM	File folder	
1429_1.csv	9/20/2019 10:16 PM	Microsoft Excel C...	47,846 KB
archive (13).zip	5/8/2025 12:12 PM	Compressed (zipp...	16,650 KB
bitext_cleaned.csv	5/12/2025 10:41 PM	Microsoft Excel C...	18,667 KB
Bitext_Sample_Customer_Support_Traini...	9/23/2023 1:55 AM	Microsoft Excel C...	18,753 KB
cleaned_amazon_reviews.csv	5/8/2025 12:22 PM	Microsoft Excel C...	7,167 KB
customer-support-llm-chatbot-training-...	5/8/2025 5:48 PM	Compressed (zipp...	2,731 KB
Datafiniti_Amazon_Consumer_Reviews_o...	9/20/2019 10:16 PM	Microsoft Excel C...	97,226 KB
Datafiniti_Amazon_Consumer_Reviews_o...	9/20/2019 10:16 PM	Microsoft Excel C...	259,418 KB
intent_classifier.pkl	5/8/2025 6:11 PM	PKL File	665 KB
Locdacbiet.py	5/9/2025 9:31 PM	JetBrains PyCharm	0 KB
trainChatbot.py	5/8/2025 6:11 PM	JetBrains PyCharm	1 KB
Tudacbiet.ipynb	5/9/2025 9:45 PM	Jupyter Source File	21 KB
Xulybitext.py	5/8/2025 6:00 PM	JetBrains PyCharm	2 KB
Xulydata.py	5/8/2025 12:21 PM	JetBrains PyCharm	2 KB

Hình 4.1.1: Thư Mục Dataset

Từ file Xulybitext.py, ta có Các intent và các Category như sau:

```
Run Xulybitext
C:\Users\LENOVO\AppData\Local\Programs\Python\Python311\python.exe "D:\2024-2025\HK2\Xu ly ngon ngu tu nhien\Week16\dataset\Xulybitext.py"
Categories: ['ORDER', 'SHIPPING', 'CANCEL', 'INVOICE', 'PAYMENT', 'REFUND', 'FEEDBACK',
'CONTACT', 'ACCOUNT', 'DELIVERY', 'SUBSCRIPTION']
Intents: ['cancel_order', 'change_order', 'change_shipping_address',
'check_cancellation_fee', 'check_invoice', 'check_payment_methods',
'check_refund_policy', 'complaint', 'contact_customer_service',
'contact_human_agent', 'create_account', 'delete_account',
'delivery_options', 'delivery_period', 'edit_account', 'get_invoice',
'get_refund', 'newsletter_subscription', 'payment_issue', 'place_order',
'recover_password', 'registration_problems', 'review',
'set_up_shipping_address', 'switch_account', 'track_order', 'track_refund']
Cleaned Bitext dataset saved to: bitext_cleaned.csv
Process finished with exit code 0
```

Hình 4.1.2: Intent Và Category

Kết quả cho thấy có 12 chủ đề (Category) lớn mà chatbot có thể xử lý bao gồm:







- Đơn hàng
- Vận chuyển
- Hủy đơn
- Hóa đơn
- Thanh toán

- Hoàn tiền
- Phản hồi
- Liên hệ
- Tài khoản
- Thời gian giao hàng
- Đăng ký nhận thông báo

Có 26 ý định (Intent) mà chatbot có thể xử lý:





- Hủy đơn hàng, thay đổi đơn hàng, đặt hàng, theo dõi đơn hàng.
- Thay đổi địa chỉ giao hàng, thiết lập địa chỉ giao hàng.
- Kiểm tra phí hủy đơn, kiểm tra hóa đơn, lấy hóa đơn.
- Kiểm tra phương thức thanh toán, lỗi thanh toán.
- Kiểm tra chính sách hoàn tiền, yêu cầu hoàn tiền, theo dõi hoàn tiền.
- Khiếu nại, liên hệ với dịch vụ khách hàng, yêu cầu nói chuyện với nhân viên.
- Tạo tài khoản, xóa tài khoản, sửa thông tin tài khoản, khôi phục mật khẩu, lỗi khi đăng ký, chuyển đổi tài khoản.
- Tùy chọn giao hàng, thời gian giao hàng.
- Đăng ký nhận bản tin.
- Đánh giá sản phẩm.

Thư mục Sentiment: Chứa file vector hóa các đoạn bình luận, các sản phẩm chất lượng cao, tỉ lệ đánh giá cùng độ tương quan của các sản phẩm:

Name	Date modified	Type	Size
 cleaned_amazon_reviews.csv	5/8/2025 12:22 PM	Microsoft Excel C...	7,167 KB
 high_quality_products.csv	5/8/2025 1:36 PM	Microsoft Excel C...	1 KB
 sentiment.py	5/8/2025 1:36 PM	JetBrains PyCharm	5 KB
 sentiment_summary.csv	5/8/2025 1:36 PM	Microsoft Excel C...	1 KB
 tfidf_matrix.pkl	5/8/2025 1:36 PM	PKL File	306 KB
 tfidf_vectorizer.pkl	5/8/2025 1:36 PM	PKL File	161 KB






**Hình 4.1.3: Thư Mục Sentiment**

Tương tự với folder Collaborative Recommendation cũng chứa ma trận tương quan giữa các sản phẩm được tính theo điểm đánh giá:

Name	Date modified	Type	Size
 cleaned_amazon_reviews.csv	5/8/2025 12:22 PM	Microsoft Excel C...	7,167 KB
 Collaborative.py	5/8/2025 12:27 PM	JetBrains PyCharm	4 KB
 product_profiles.pkl	5/8/2025 12:27 PM	PKL File	3 KB
 similarity_matrix.pkl	5/8/2025 12:27 PM	PKL File	5 KB

**Hình 4.1.4: Thư Mục Collaborative Recommendation**

Đối với thư mục Sentiment\_User: Ta có mô hình phân loại cảm xúc cho người dùng:

Name	Date modified	Type	Size
 bitext_cleaned.csv	5/8/2025 6:00 PM	Microsoft Excel C...	18,667 KB
 sentiment_model_user.pkl	5/9/2025 10:28 PM	PKL File	58 KB
 sentiment_user.ipynb	5/9/2025 10:27 PM	Jupyter Source File	64 KB
 sentiment_user.py	5/9/2025 10:18 PM	JetBrains PyCharm	2 KB
 tfidf_vectorizer_user.pkl	5/9/2025 10:28 PM	PKL File	79 KB

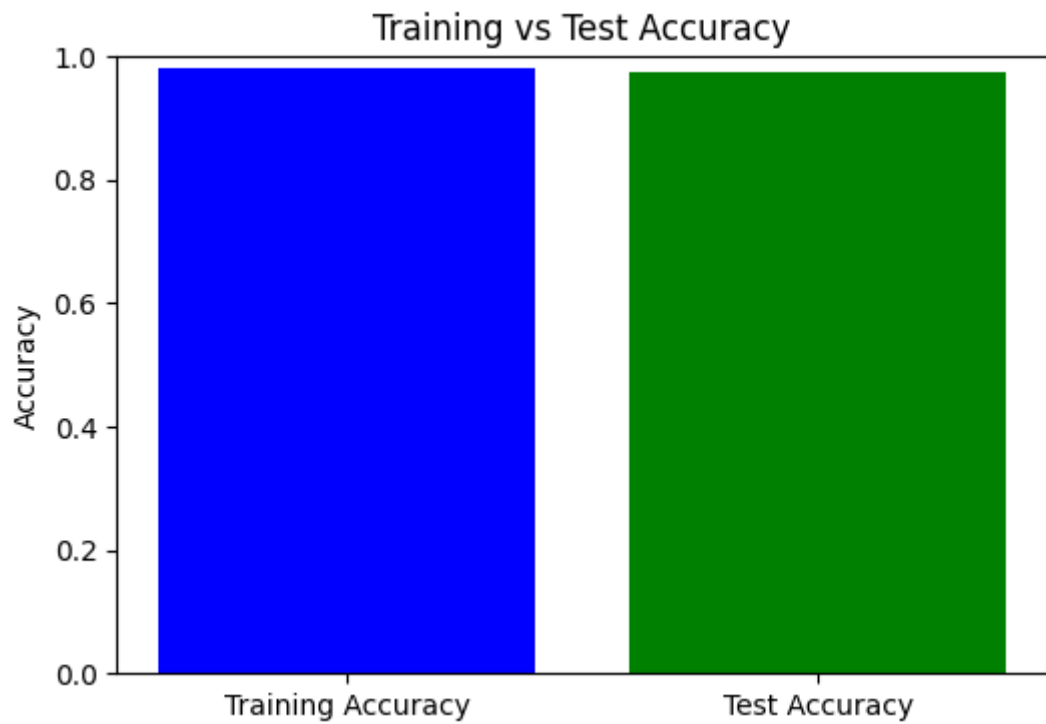
**Hình 4.1.5: Thư Mục Sentiment\_User**

Kèm theo đó là báo cáo phân loại cảm xúc của người dùng được huấn luyện từ tập bitext:

Model Evaluation:				
	precision	recall	f1-score	support
negative	0.93	0.95	0.94	783
neutral	0.99	0.98	0.99	2937
positive	0.97	0.98	0.97	1655
accuracy			0.98	5375
macro avg	0.96	0.97	0.97	5375
weighted avg	0.98	0.98	0.98	5375

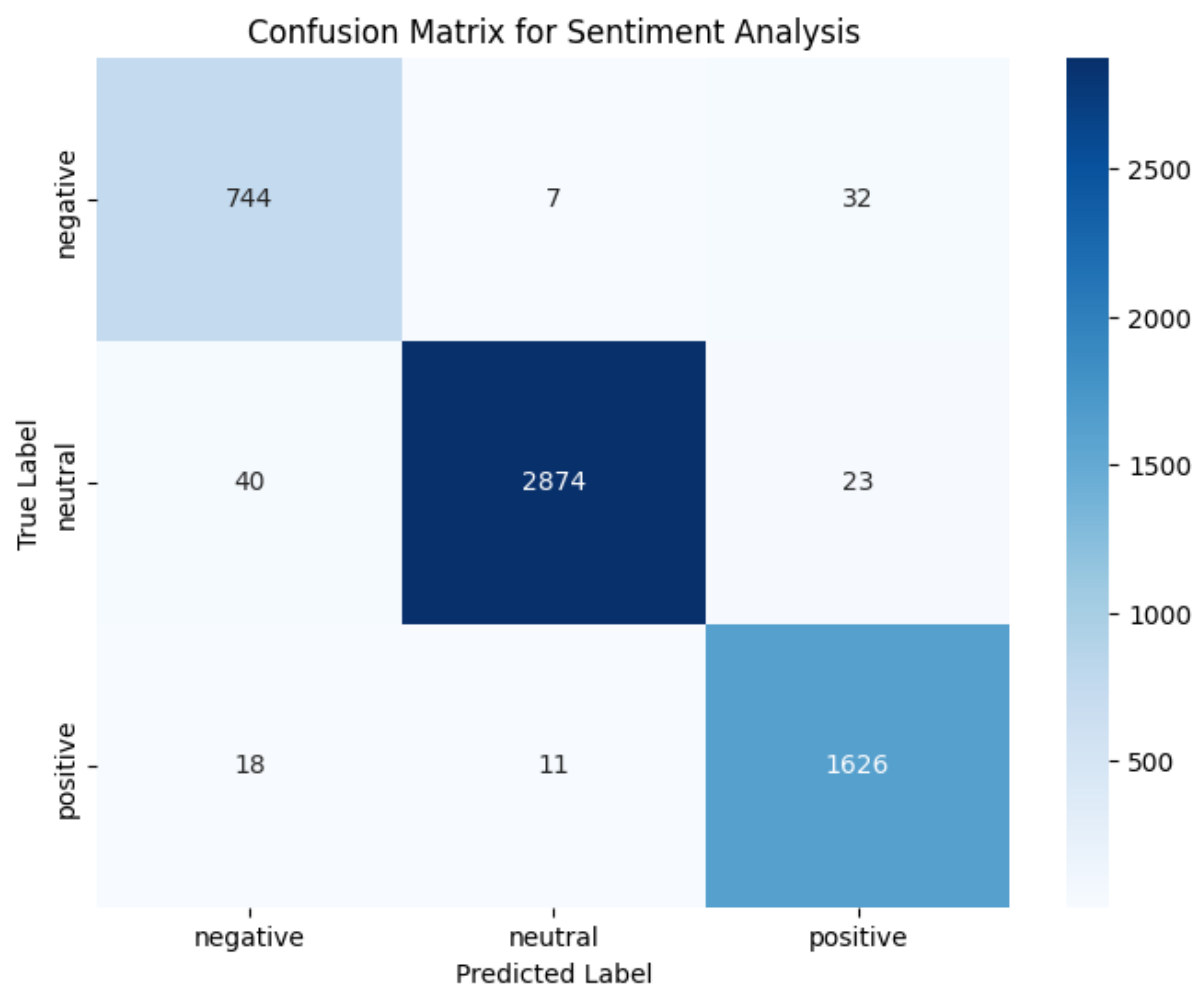
Model and vectorizer saved successfully.  
 Training Accuracy: 0.9810  
 Test Accuracy: 0.9756

**Hình 4.1.6: Báo Cáo Phân Loại**



*Hình 4.1.7: Accuracy Của Cột Training Và Cột Test*

Dựa trên báo cáo phân loại cùng với Biểu đồ độ chính xác giữa tập huấn luyện và tập kiểm tra có độ chính xác cao, nghĩa là mô hình hoạt động rất tốt với tập dữ liệu `bitext_cleaned.csv`. Và ta có ma trận nhầm lẫn:



*Hình 4.1.8: Ma Trận Nhầm Lẫn*

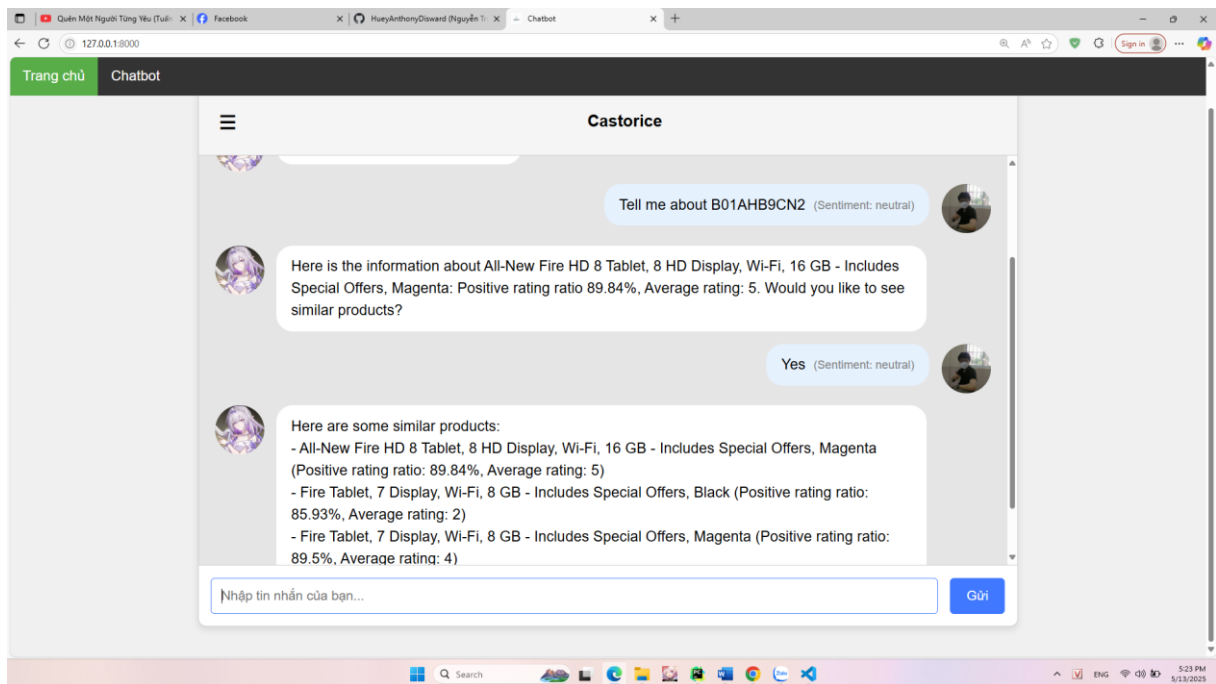
## 4.2. Giao diện

Cuối cùng là giao diện chatbot xây dựng bằng Django:

Giao diện lúc bắt đầu:

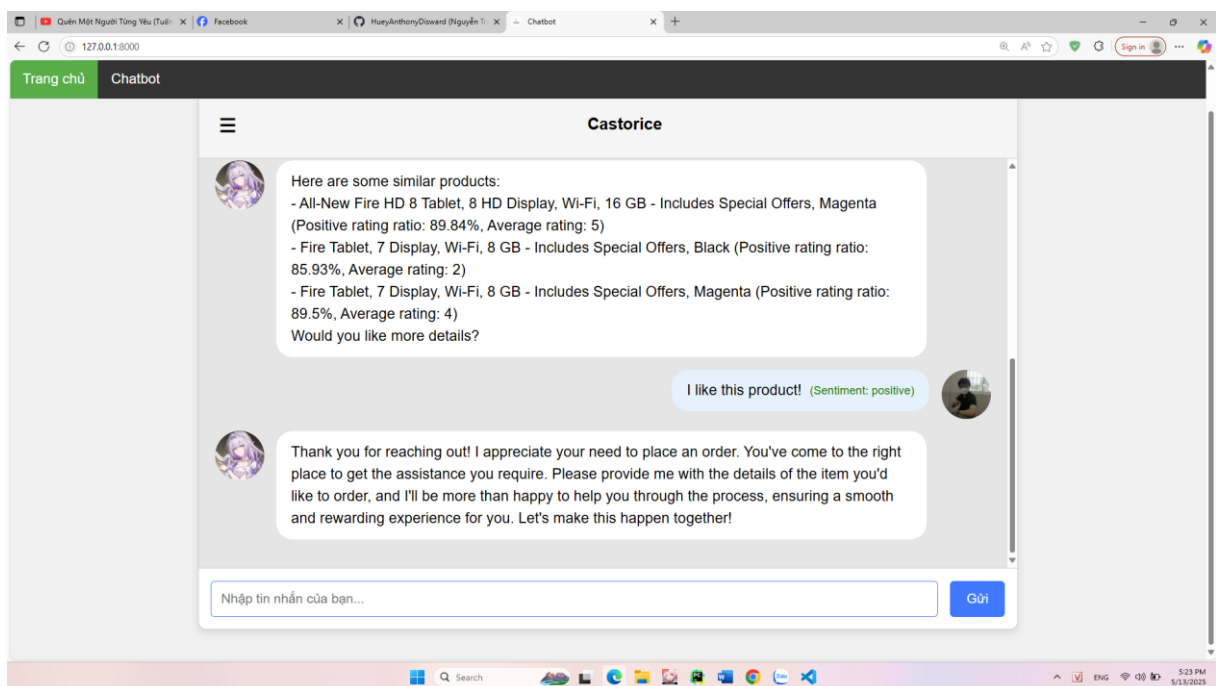
*Hình 4.2.1: Giao Diện Bắt Đầu*

Tư vấn thông tin và gợi ý sản phẩm:



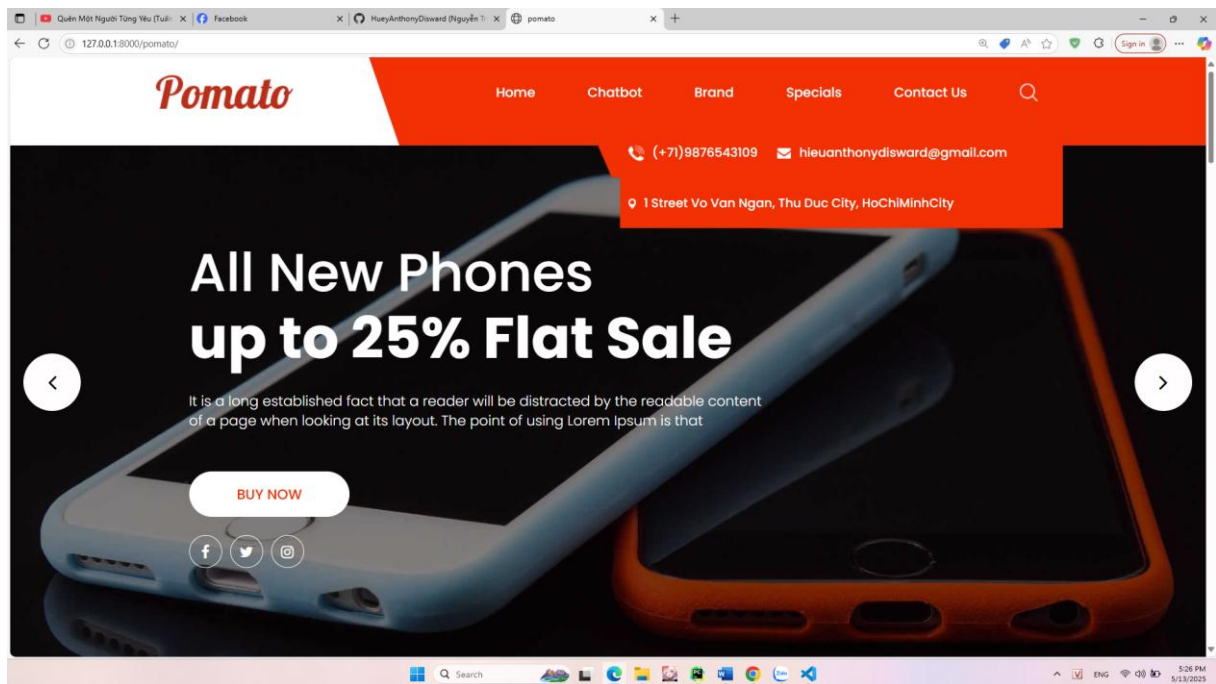
*Hình 4.2.2: Tư Vấn Và Gợi Ý Sản Phẩm*

Nhận diện cảm xúc tích cực:



*Hình 4.2.3: Nhận Diện Cảm Xúc Tích Cực*

Trang chủ của web:



*Hình 4.2.4: Trang Chủ Web*

## PHẦN 5: KẾT LUẬN

### 5.1. Kết luận

Dự án phát triển hệ thống giao tiếp tự động tích hợp khả năng nhận diện cảm xúc và khuyến nghị dựa trên mặt hàng đã khẳng định vai trò then chốt của công nghệ thông minh trong việc nâng cao chất lượng tương tác số. Kết hợp các kỹ thuật tiên tiến như xử lý ngôn ngữ tự nhiên, học máy và phân tích dữ liệu, hệ thống không chỉ đáp ứng nhu cầu về giao tiếp nhạy bén và cá nhân hóa mà còn định hình lại cách doanh nghiệp xây dựng mối liên kết với khách hàng, đặc biệt trong thương mại điện tử và dịch vụ hỗ trợ. Giải pháp này thể hiện sự đột phá công nghệ, đồng thời làm nổi bật giá trị của sự kết hợp giữa trí tuệ nhân tạo và yếu tố con người, mang lại tác động sâu rộng trong kỷ nguyên số.

Ở khía cạnh học thuật, dự án đã đóng góp một nghiên cứu quan trọng vào lĩnh vực giao tiếp người-máy, làm sáng tỏ tiềm năng của việc tích hợp các phương pháp phân tích cảm xúc và khuyến nghị cá nhân hóa trong một hệ thống thống nhất. Việc ứng dụng các công cụ như VADER, Logistic Regression và Cosine Similarity không chỉ làm giàu kiến thức về xử lý ngôn ngữ tự nhiên mà còn cung cấp một trường hợp thực tiễn giá trị cho cộng đồng nghiên cứu. Kết quả đạt được đã chứng minh hiệu quả của các phương pháp liên ngành, củng cố vai trò của trí tuệ nhân tạo trong việc phát triển các công cụ giao tiếp thông minh và tạo nguồn tài liệu tham khảo quý báu cho các nghiên cứu tương lai.

Hệ thống này mang lại lợi ích mang tính thực tiễn cho doanh nghiệp và người dùng cuối. Với các tổ chức hoạt động trong môi trường cạnh tranh, công cụ này tối ưu hóa việc xử lý yêu cầu, nâng cao sự hài lòng của khách hàng và giảm chi phí thông qua tự động hóa hiệu quả. Người dùng, từ khách hàng mua sắm trực tuyến đến những người sử dụng dịch vụ hỗ trợ, được trải nghiệm các tương tác được cá nhân hóa, tạo cảm giác được quan tâm và hỗ trợ kịp thời. Hệ thống không chỉ cải thiện trải nghiệm người dùng mà còn thiết lập một chuẩn mực mới cho các dịch vụ số, kết hợp tính chính xác và sự gần gũi. Trong bối cảnh số hóa ngày càng mở rộng, dự án đã khẳng định giá trị như một giải pháp thực tiễn, đáp ứng nhu cầu kết nối hiệu quả và tạo ra tác động tích cực đến cả cộng đồng doanh nghiệp lẫn người tiêu dùng, nhấn mạnh vai trò của công nghệ trong việc giải quyết các thách thức thực tế.



## 5.2. Hướng phát triển

Hệ thống giao tiếp tự động tích hợp nhận diện cảm xúc và khuyến nghị cá nhân hóa mở ra nhiều triển vọng để nâng cao tính linh hoạt và giá trị ứng dụng trong bối cảnh công nghệ số không ngừng tiến hóa. Một định hướng quan trọng là mở rộng các phương thức tương tác, vượt ra ngoài văn bản bằng cách tích hợp khả năng xử lý giọng nói hoặc nhận diện hình ảnh, nhằm mang lại trải nghiệm giao tiếp tự nhiên và tiện lợi hơn. Đồng thời, việc nâng cấp các thuật toán phân tích để xử lý các ngữ cảnh giao tiếp phức tạp hơn, như nhận diện các sắc thái cảm xúc tinh vi hoặc phân tích các cuộc đối thoại dài, sẽ giúp hệ thống phản hồi chính xác và cá nhân hóa hơn, đáp ứng tốt các kịch bản thực tế đa dạng.

Bên cạnh đó, hệ thống có tiềm năng được ứng dụng vào các lĩnh vực mới ngoài thương mại điện tử và dịch vụ khách hàng, mở rộng phạm vi tác động. Trong giáo dục, công cụ này có thể được tinh chỉnh để hỗ trợ học tập cá nhân hóa, nhận biết trạng thái tâm lý của học viên và đề xuất tài liệu hoặc phương pháp học phù hợp, từ đó tối ưu hóa kết quả học tập. Trong y tế, hệ thống có thể đóng vai trò tư vấn sức khỏe tâm lý, cung cấp các phản hồi đồng cảm và giải pháp dựa trên cảm xúc của bệnh nhân, góp phần nâng cao chất lượng chăm sóc. Hơn nữa, trong quản lý nhân sự, doanh nghiệp có thể tận dụng hệ thống để cải thiện giao tiếp nội bộ, như đánh giá tâm trạng nhân viên và gợi ý các chương trình đào tạo hoặc phúc lợi phù hợp. Những ứng dụng này không chỉ khẳng định tính linh hoạt của hệ thống mà còn đáp ứng nhu cầu chuyên biệt của từng ngành.

Việc tối ưu hóa hiệu suất và tích hợp các tiến bộ mới sẽ đảm bảo hệ thống đáp ứng các yêu cầu ngày càng cao. Cải tiến thuật toán để nâng cao độ chính xác và tốc độ xử lý, đặc biệt với dữ liệu đa ngôn ngữ hoặc khối lượng lớn, sẽ giúp hệ thống hoạt động hiệu quả trong môi trường toàn cầu. Việc khai thác các nguồn dữ liệu đa dạng, như phản hồi từ mạng xã hội hoặc dữ liệu thời gian thực, có thể tăng cường khả năng hiểu ngữ cảnh và đưa ra khuyến nghị phù hợp. Ngoài ra, tích hợp hệ thống với các công nghệ tiên tiến, như thiết bị IoT hoặc giao diện thực tế tăng cường, sẽ tạo ra các trải nghiệm tương tác phong phú và sáng tạo. Những định hướng này không chỉ củng cố giá trị của hệ thống mà còn đảm bảo khả năng thích nghi với các xu hướng công nghệ mới, mang lại lợi ích bền vững cho cả cộng đồng nghiên cứu và các bên ứng dụng thực tiễn.

## TÀI LIỆU THAM KHẢO

- [1] H. Design, "Pomato Free CSS Template," 2020. [Online]. Available: <https://www.free-css.com/free-css-templates/page262/pomato> .
- [2] A. W. Services, "What is a Chatbot?," [Online]. Available: <https://aws.amazon.com/vi/what-is/chatbot/> .
- [3] GeeksforGeeks, "Django Tutorial," [Online]. Available: <https://www.geeksforgeeks.org/django-tutorial/> .
- [4] A. Vidhya, "Introduction to Collaborative Filtering," 2024. [Online]. Available: <https://www.analyticsvidhya.com/blog/2022/02/introduction-to-collaborative-filtering/> .