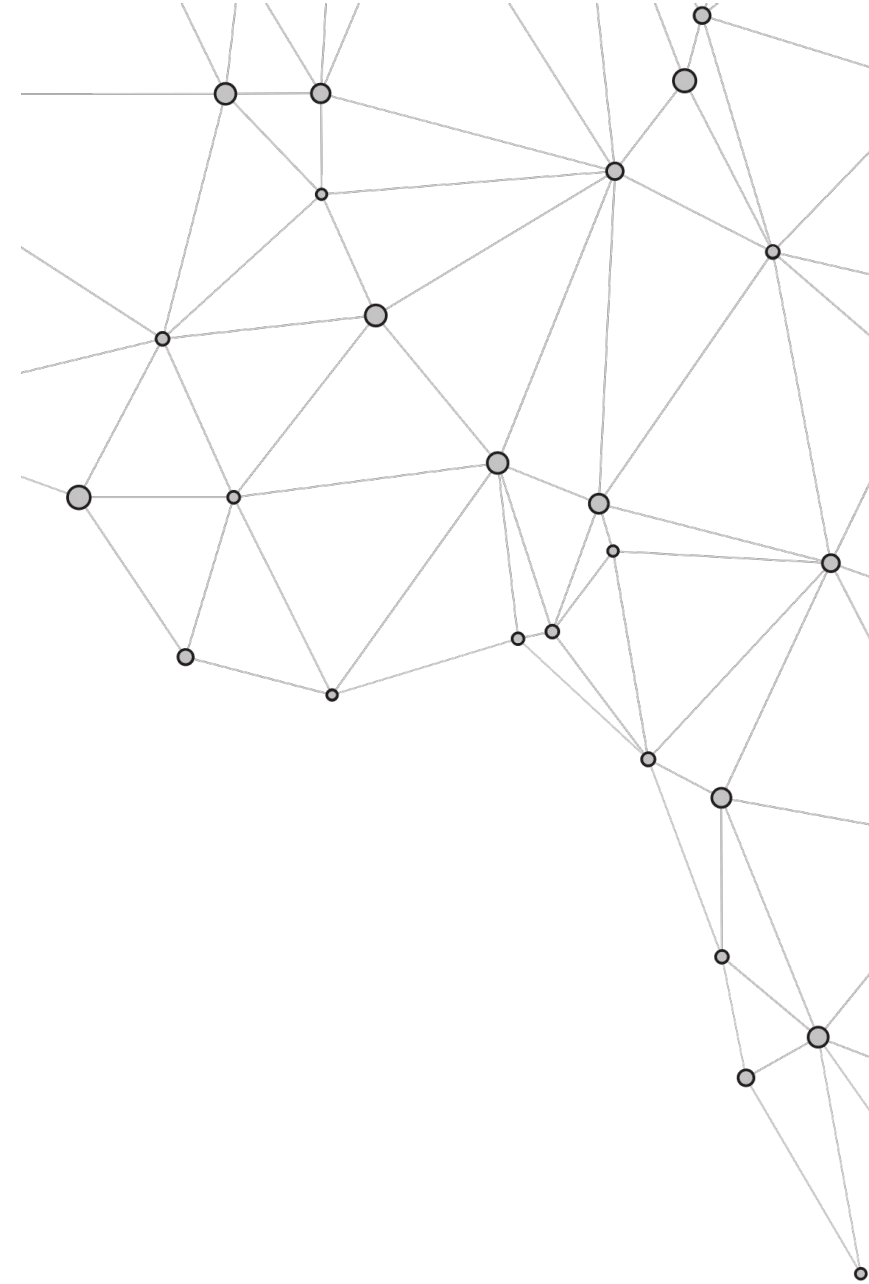


Topic Modelling

Analyze Facebook Post Topic



Danh sách thành viên

Nhóm 11

Họ và Tên: Lại Tuấn Anh

MSSV: 18020117

Email: 18020117@vnu.edu.vn

Họ và Tên: Triệu Vũ Hải

MSSV: 18020442

Email: 18020442@vnu.edu.vn

Họ và Tên: Nguyễn Thế Quân

MSSV: 18021030

Email: 18021030@vnu.edu.vn

Họ và Tên: Nguyễn Hữu Huy

MSSV: 18020644

Email: 18020644@vnu.edu.vn

Mục lục



Giới thiệu bài toán



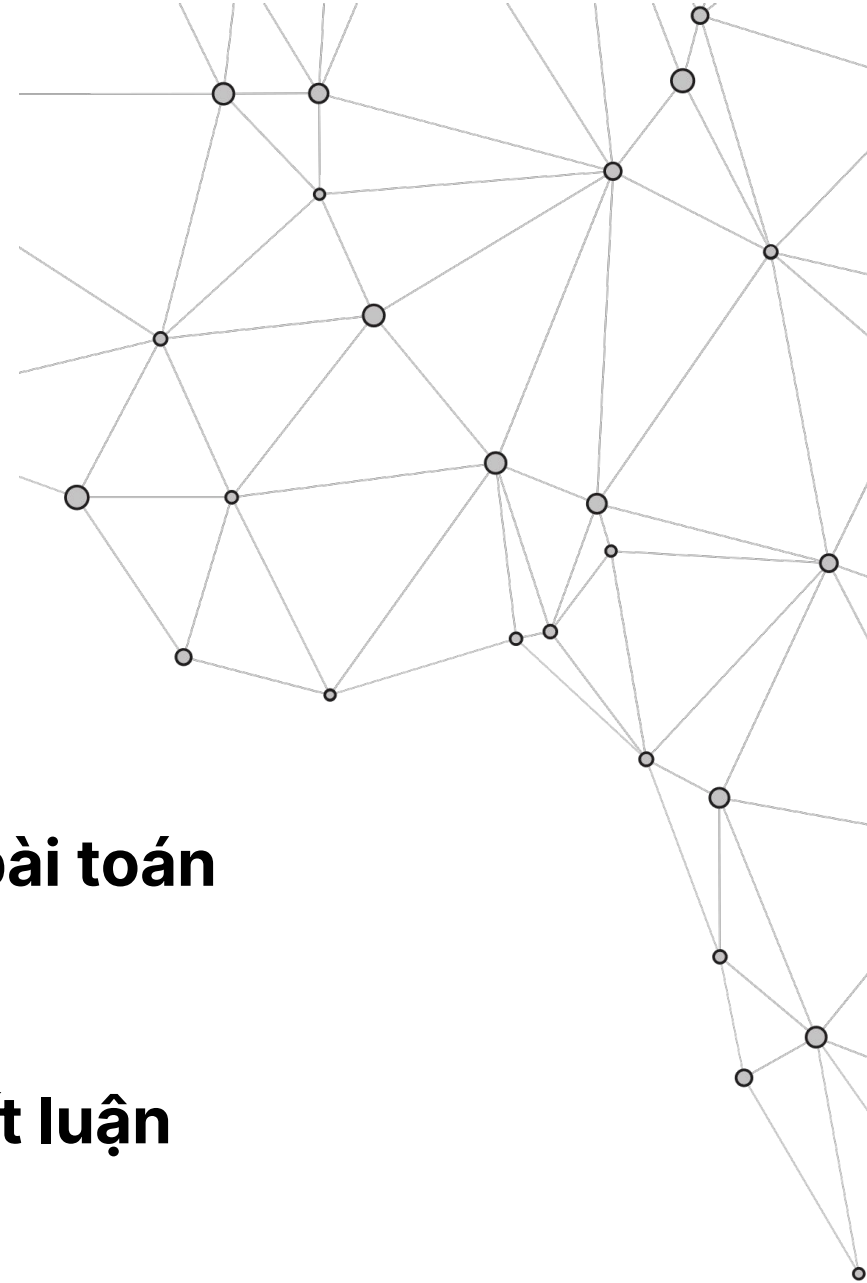
Phân tích dữ liệu



Phương pháp giải quyết bài toán

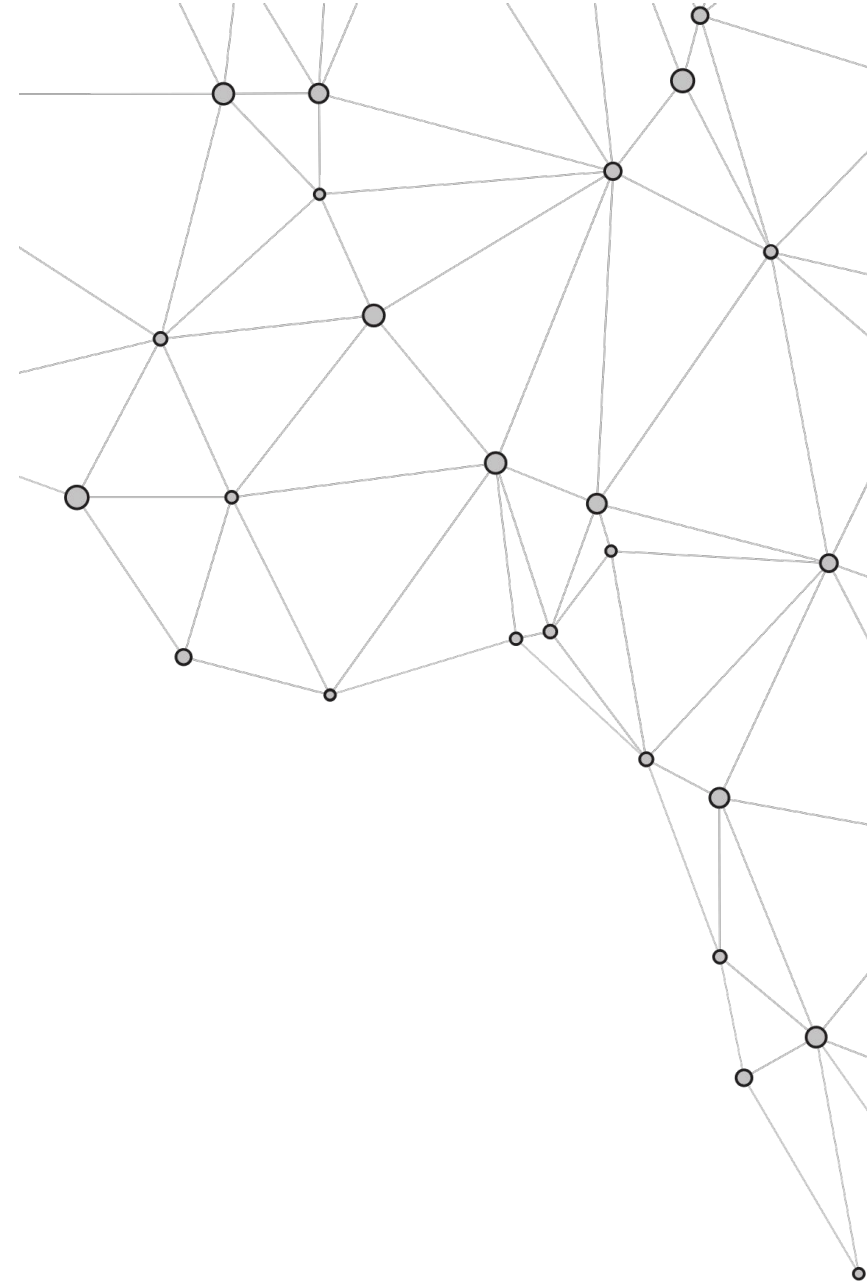


Kết quả thử nghiệm và kết luận



Giới thiệu bài toán

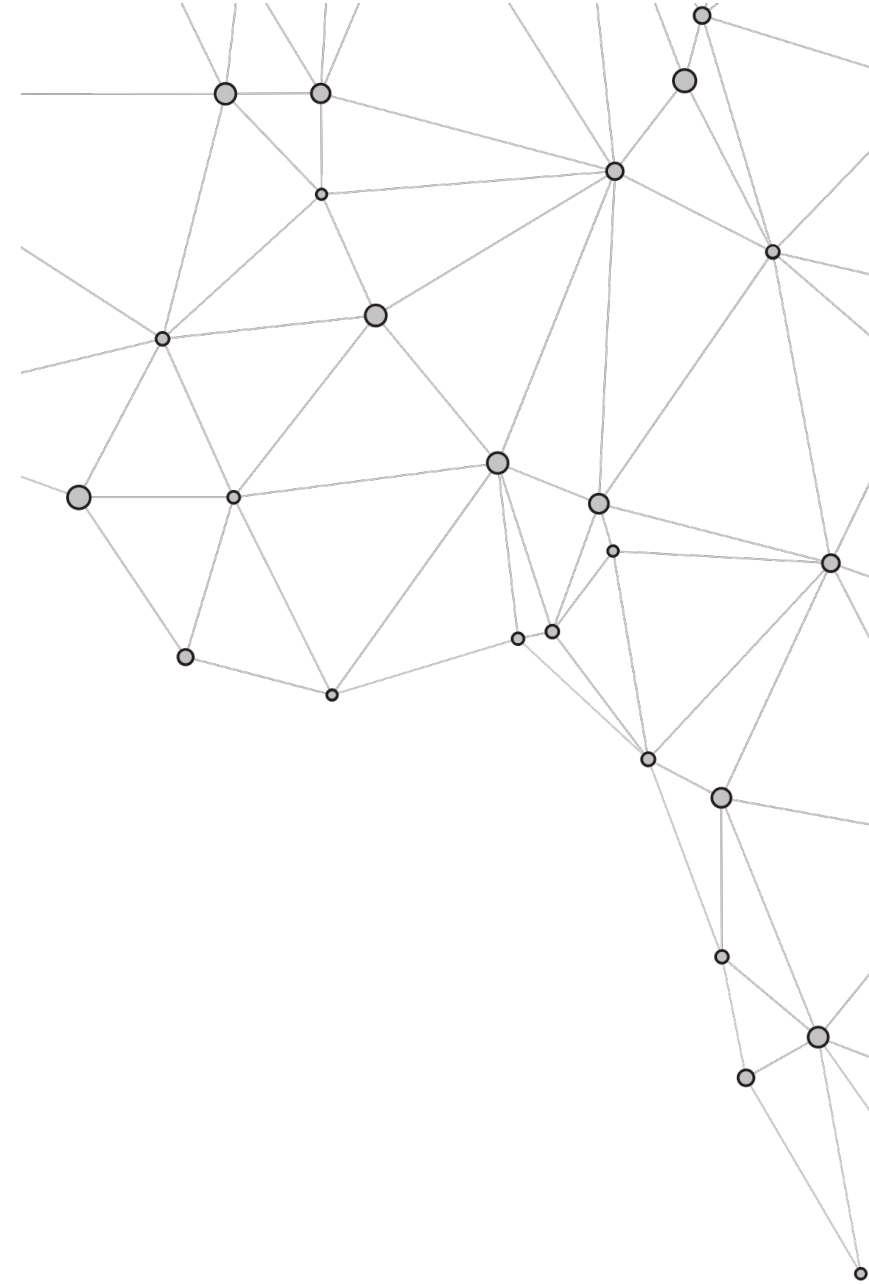
Thực hiện đánh nhãn cho tập dữ liệu văn bản bao gồm 16000 bản ghi và 23 nhãn khác nhau.



Phân tích dữ liệu

Các bước phân tích dữ liệu

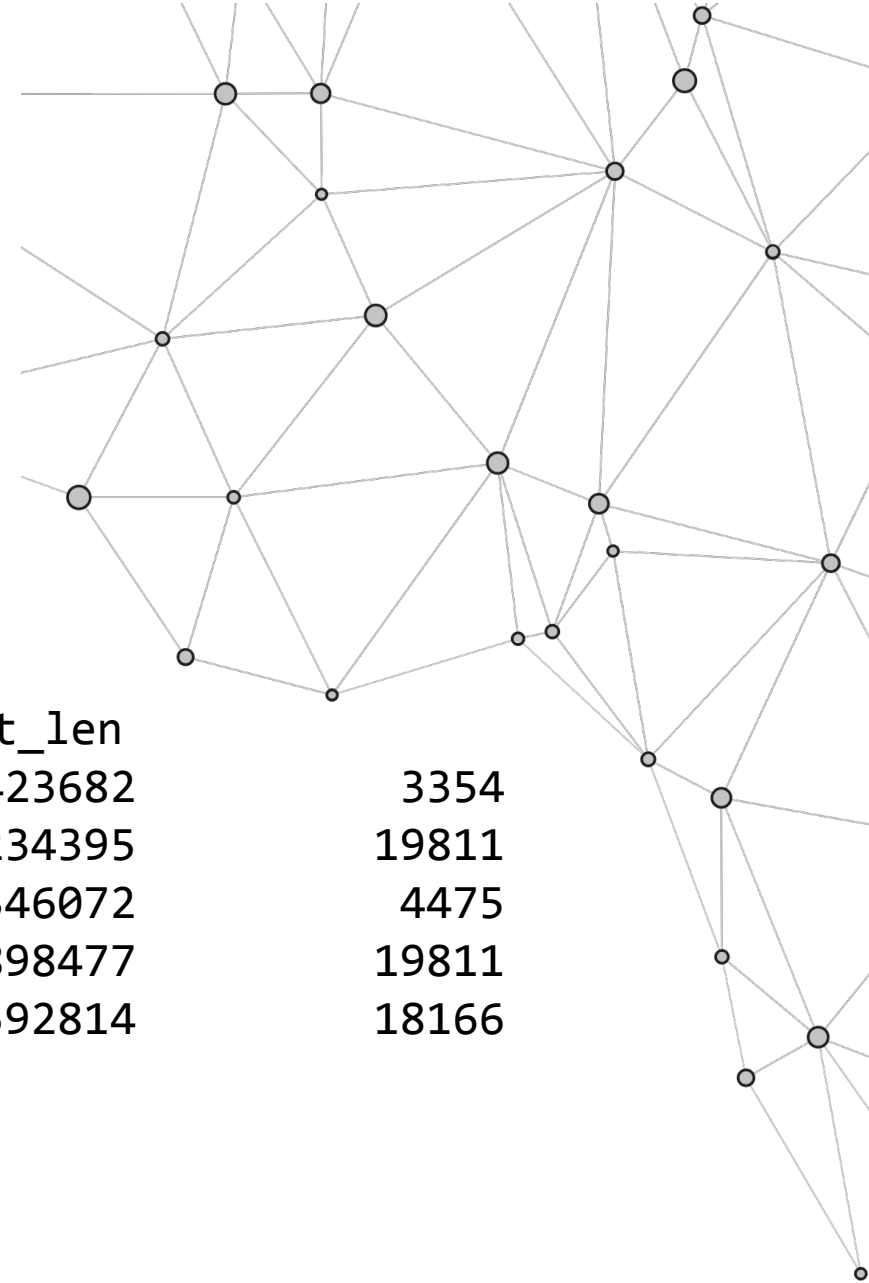
- Làm sạch dữ liệu
- Thu thập thông tin dữ liệu
- Visualize dữ liệu



Thu thập thông tin

Kết quả

	label	count	mean_text_len	max_text_len
15	__label__Nha_dat	2542	369.423682	3354
9	__label__Kinh_doanh_va_Cong_nghiep	2355	981.234395	19811
3	__label__Do_an_va_do_uong	2355	612.346072	4475
20	__label__Tai_chinh	1379	1065.898477	19811
13	__label__Mua_sam	1169	507.592814	18166
...				



Làm sạch dữ liệu

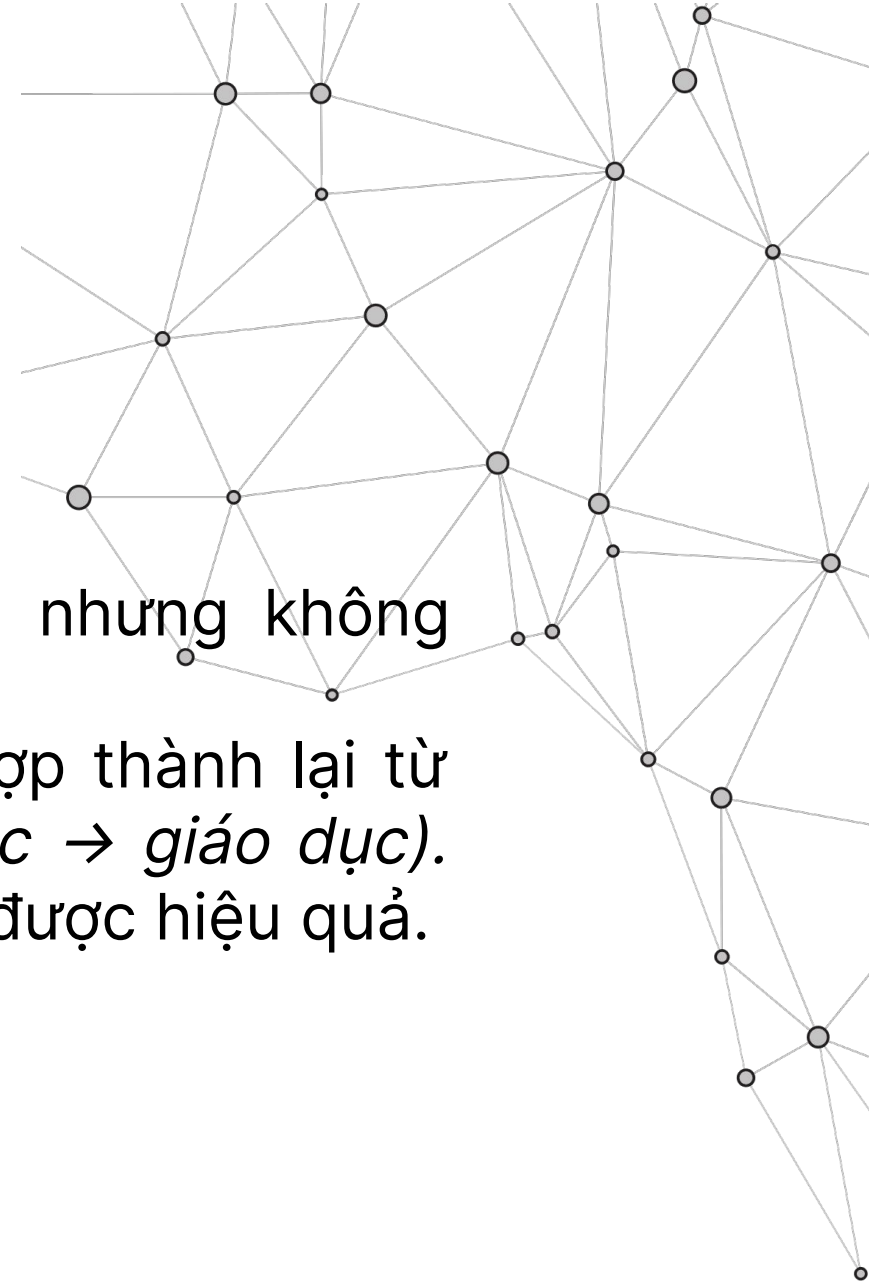
Vấn đề

- Dữ liệu là những đoạn văn bản được người dùng tạo ra. Vì vậy, dữ liệu có cấu trúc rời rạc, không tối ưu cho quá trình phân tích dữ liệu.
- Dữ liệu văn bản có nhiều kí tự gây nhiễu (emoji, số điện thoại, đường dẫn), làm ảnh hưởng đến chất lượng mô hình nếu không được xử lý.

Làm sạch dữ liệu

Vấn đề

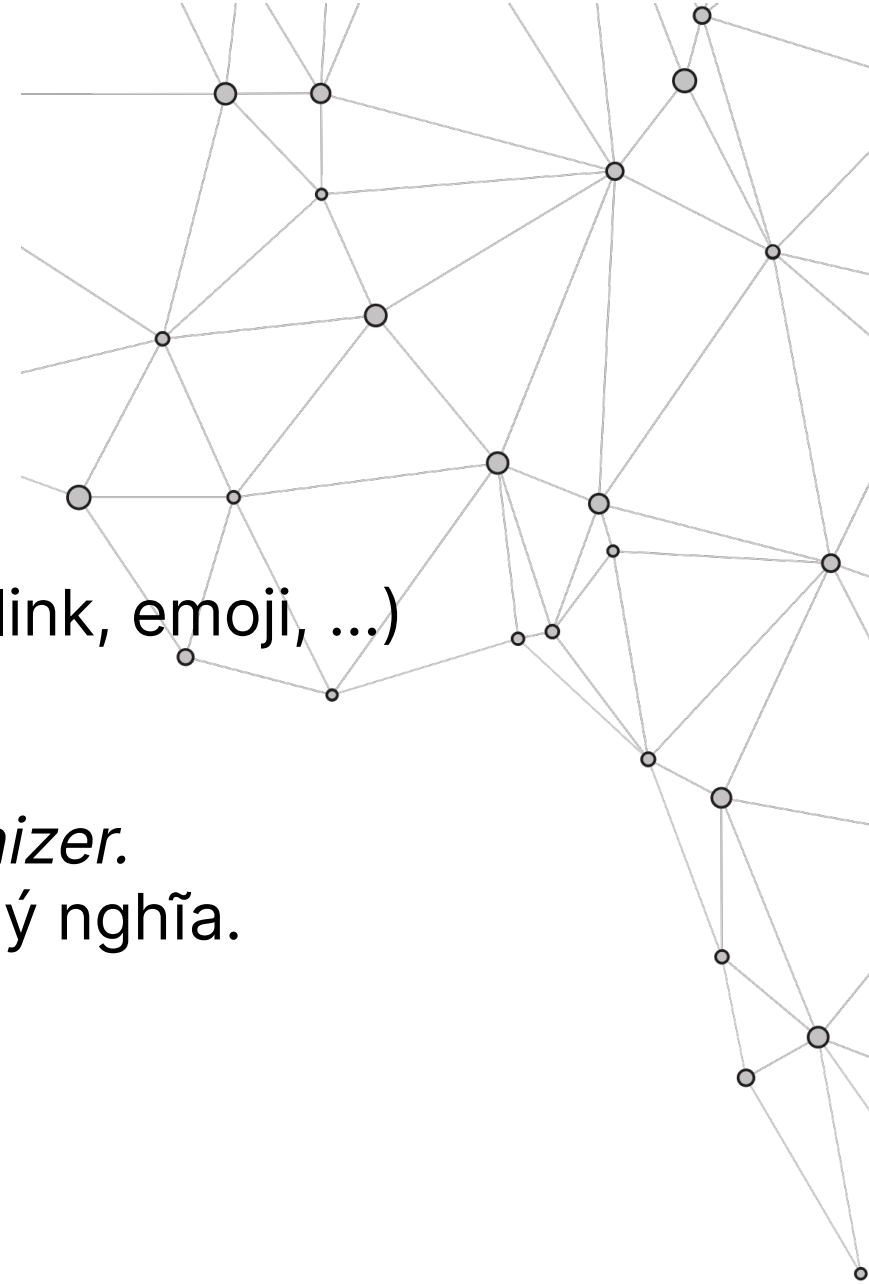
- Dữ liệu có nhiều từ được sử dụng nhiều lần nhưng không mang nhiều ý nghĩa (*thế, nhưng, tôi, mà,...*).
- Tiếng Việt có nhiều từ ghép. Từ ghép được hợp thành lại từ hai từ đơn có nghĩa hoặc vô nghĩa (*giáo + dục → giáo dục*). Do vậy, ta cần nối các từ ghép lại để phân tích được hiệu quả.



Làm sạch dữ liệu

Hướng giải quyết

- Xóa dữ liệu nhiễu bằng cách sử dụng RegExp. (link, emoji, ...)
- Đưa dữ liệu về chữ in thường.
- Chuẩn hóa dấu trong tiếng Việt.
- Nối hai từ đơn thành một từ ghép bằng *ViTokenizer*.
- Xóa các từ xuất hiện nhiều lần nhưng không có ý nghĩa.



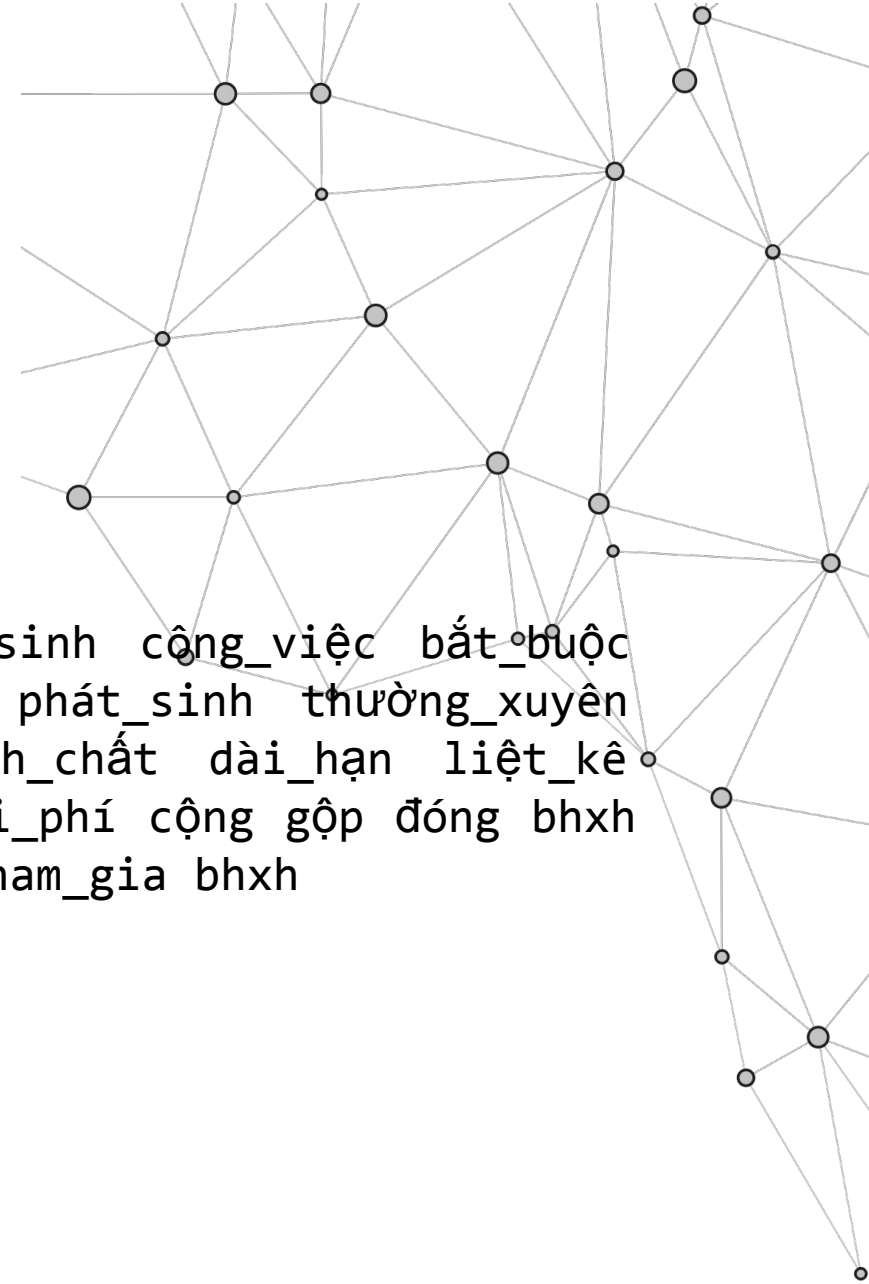
Làm sạch dữ liệu

Kết quả

Cả nhà cho em hỏi với ạ 🥰!! Công ty em hay phát sinh một số công việc bắt buộc phải cho nhân viên tăng ca, phát sinh không thường xuyên nhưng công việc mang tính chất dài hạn. Em liệt kê chi phí tăng ca đó vào loại chi phí nào để không cộng gộp đóng BHXH hoặc phải ký hợp đồng nào để không phải tham gia BHXH ạ? <3

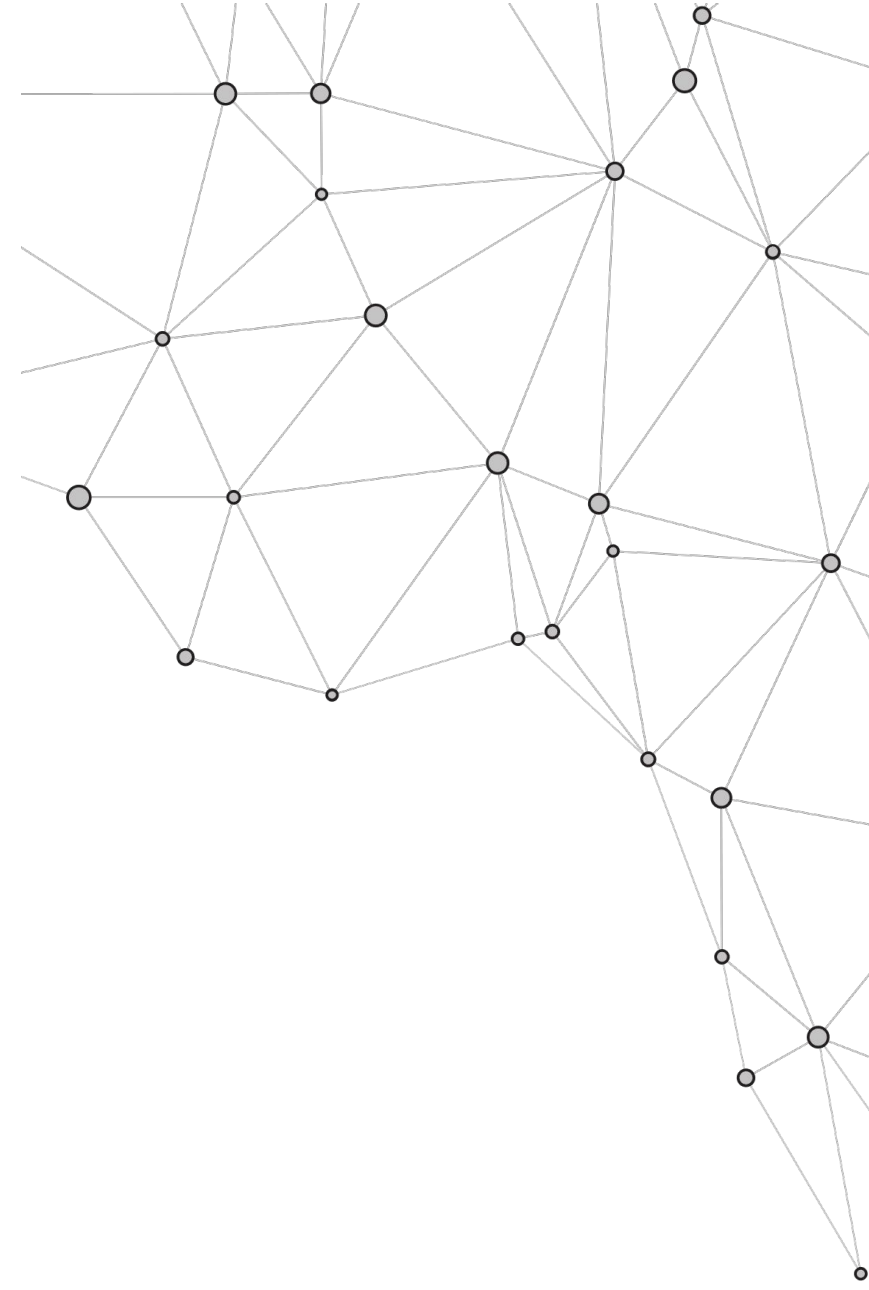
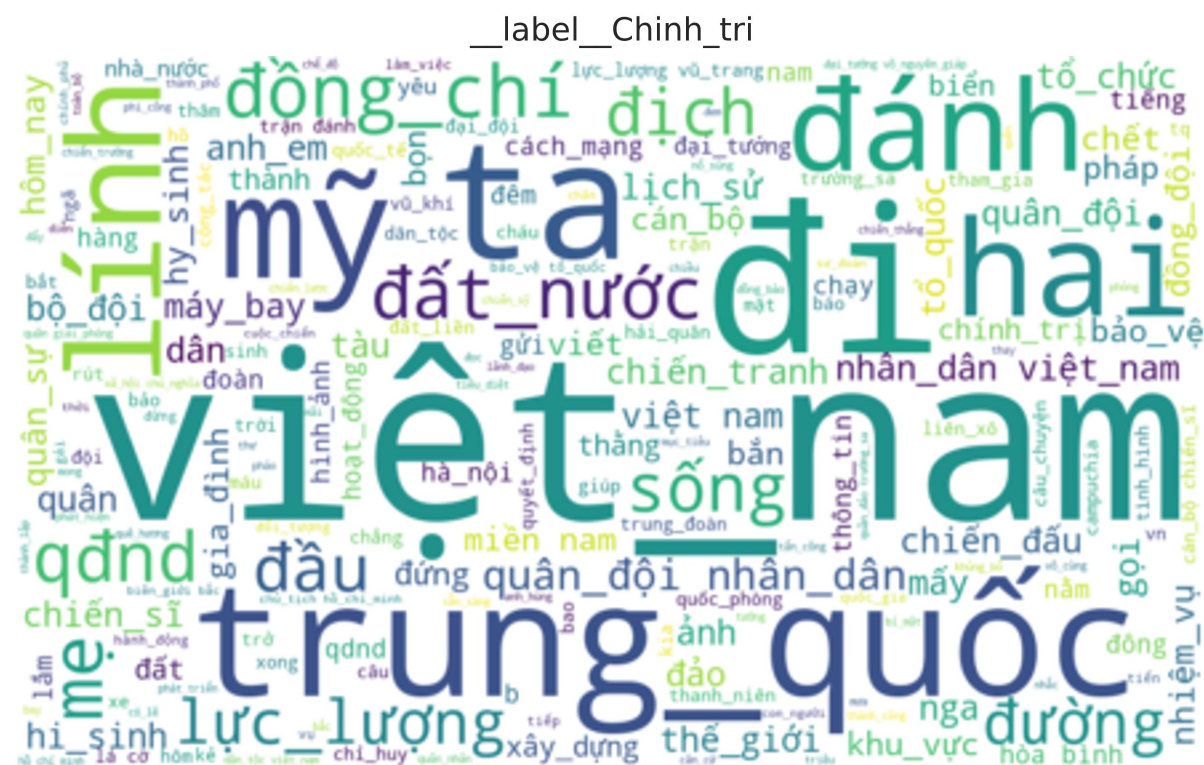


công_ty phát_sinh công_việc bắt_buộc
nhân_viên ca phát_sinh thường_xuyên
công_việc tính_chất dài_hạn liệt_kê
chi_phí ca chi_phí cộng gộp đóng BHXH
ký_hợp_đồng tham_gia BHXH



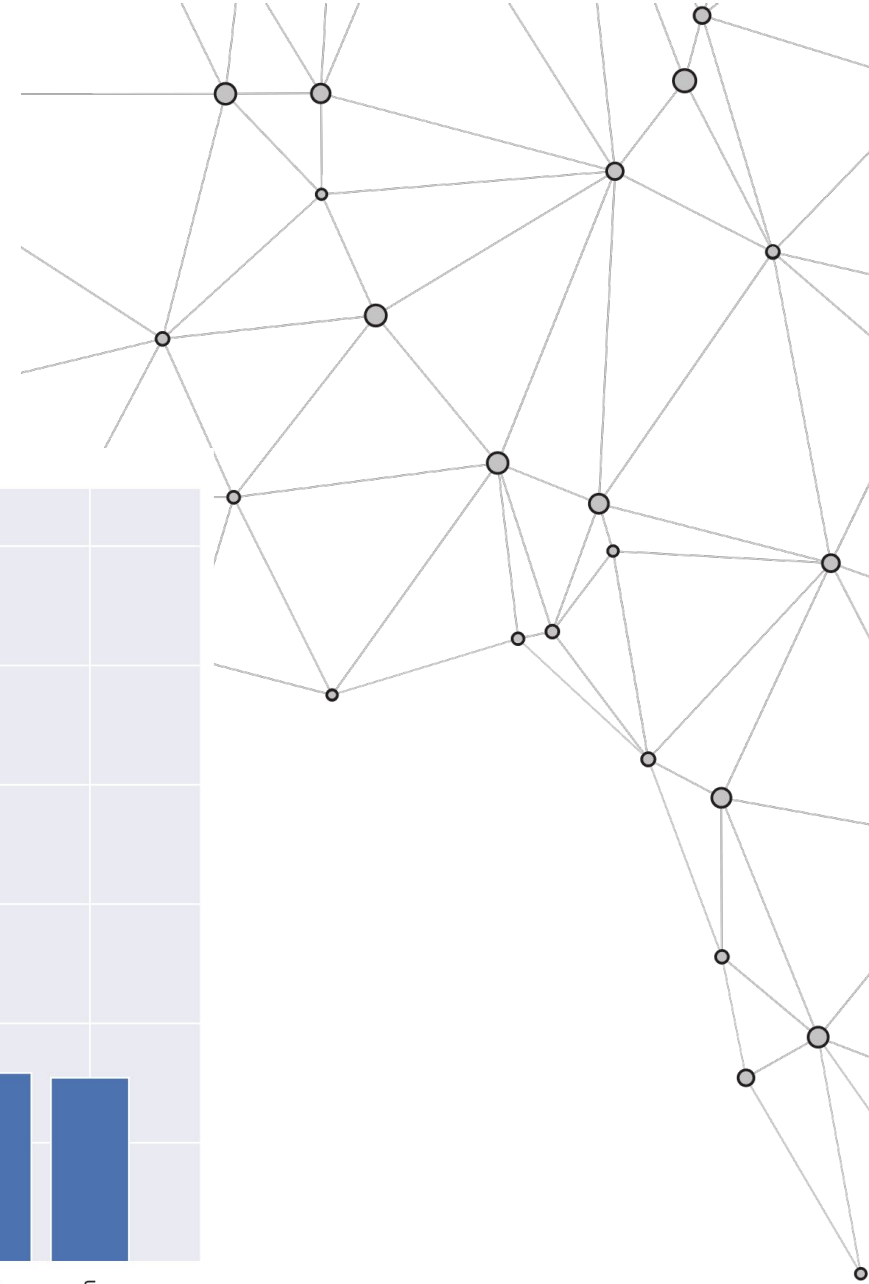
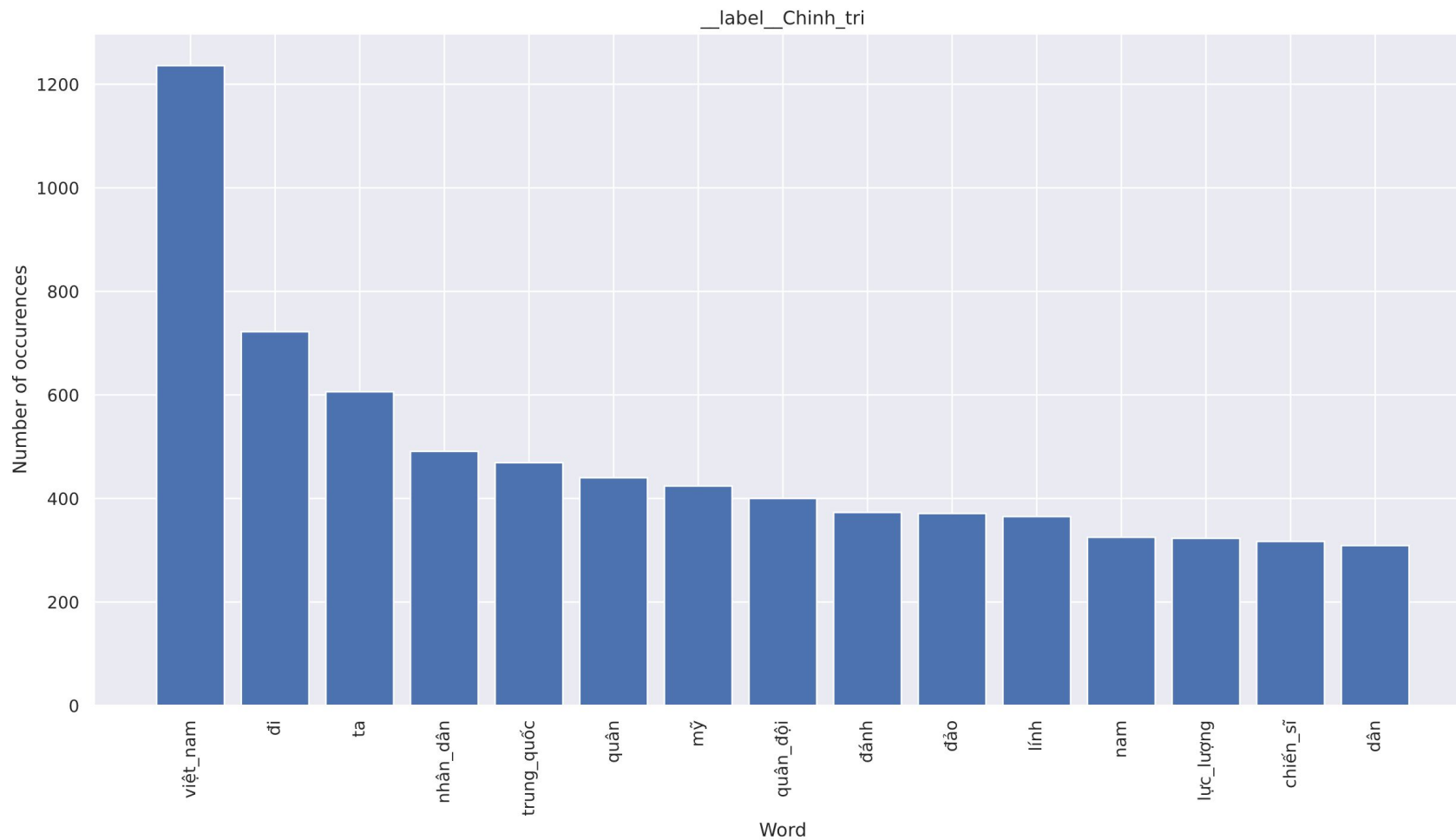
Visualize dữ liệu

Word Cloud



Visualize dữ liệu

Bar Chart



Phương pháp giải quyết bài toán

Các bước giải quyết bài toán

- Chia tập dữ liệu đã được làm sạch thành 2 tập: Train và Test theo tỉ lệ 75-25.
- Vector hóa dữ liệu đã được làm sạch bằng công thức TF-IDF.
- Từ những vector được sinh ra, huấn luyện mô hình bằng thuật toán SVM.
- Kiểm tra độ chính xác của thuật toán bằng cách so sánh kết quả dự đoán với kết quả thật.

Phân chia dữ liệu **train** và **test**

- Dựa trên dữ liệu được cung cấp, nhóm đã chia tập dữ liệu thành 2 phần:
 - **Tập train:** sử dụng để xây dựng mô hình giải quyết bài toán dựa trên TFIDF và SVM.
 - **Tập test:** dữ liệu dùng để đánh giá độ hiệu quả của mô hình, mức độ chính xác của việc phân loại dữ liệu. Dữ liệu test sẽ được bỏ nhãn).
- **Cách chia dữ liệu:** dữ liệu được chia theo tỉ lệ train:test là 75:25 với mỗi nhãn có trong tập dữ liệu để đảm bảo độ tương đồng tỉ lệ các nhãn giữa tập train và tập test.



Thuật toán TF-IDF

- **Mục đích:** Phản ánh độ quan trọng của mỗi từ hoặc n-gram đối với văn bản trên toàn bộ tài liệu đầu vào
- **TF:** Term Frequency (Tần suất xuất hiện của từ) là số lần từ xuất hiện trong văn bản

$$\text{idf}(t, D) = \log \frac{|D|}{|\{d \in D : t \in d\}|}$$

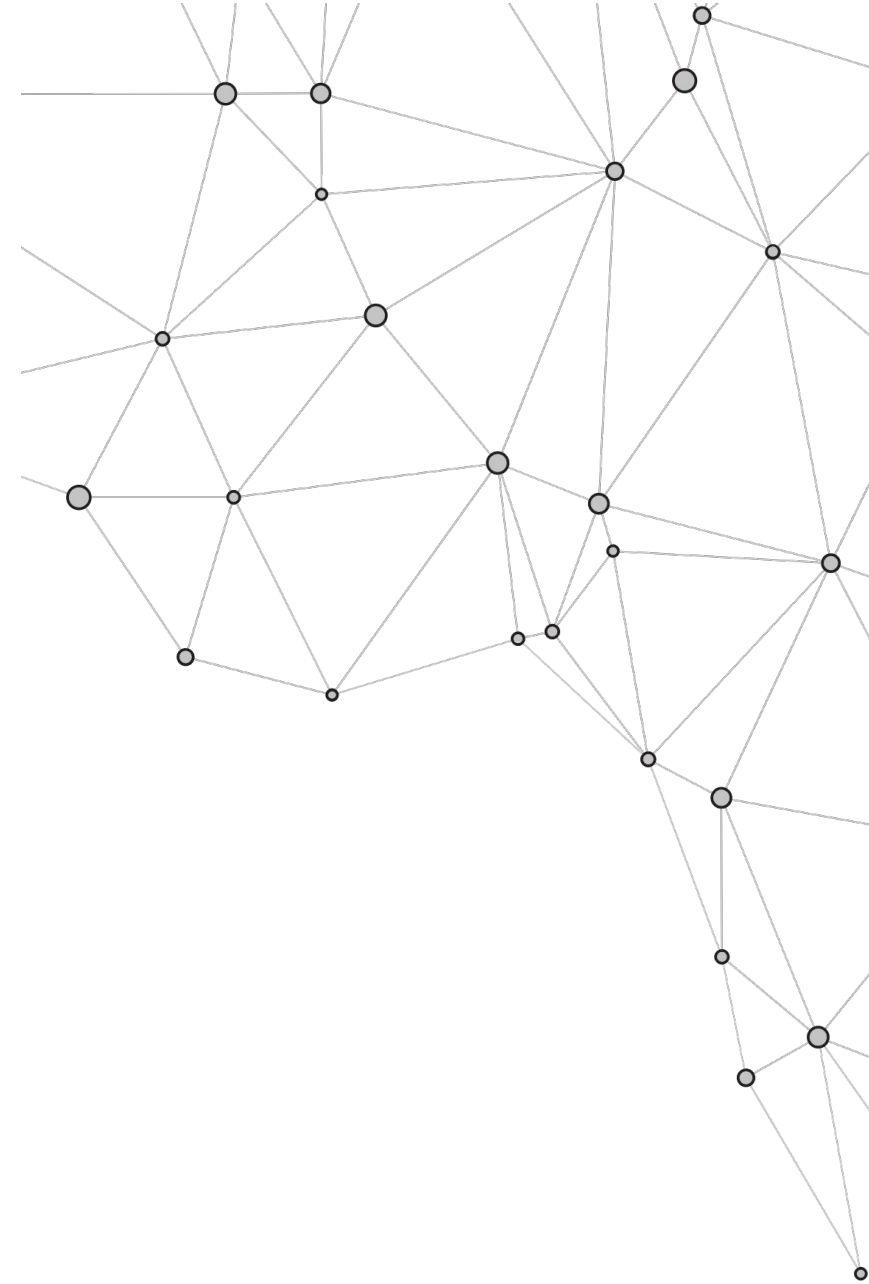
- **IDF:** Inverse Document Frequency (Nghịch đảo tần suất của văn bản), giúp đánh giá tầm quan trọng của một từ

$$\text{tf}(t, d) = \frac{f(t, d)}{\max\{f(w, d) : w \in d\}}$$

Thuật toán TF-IDF

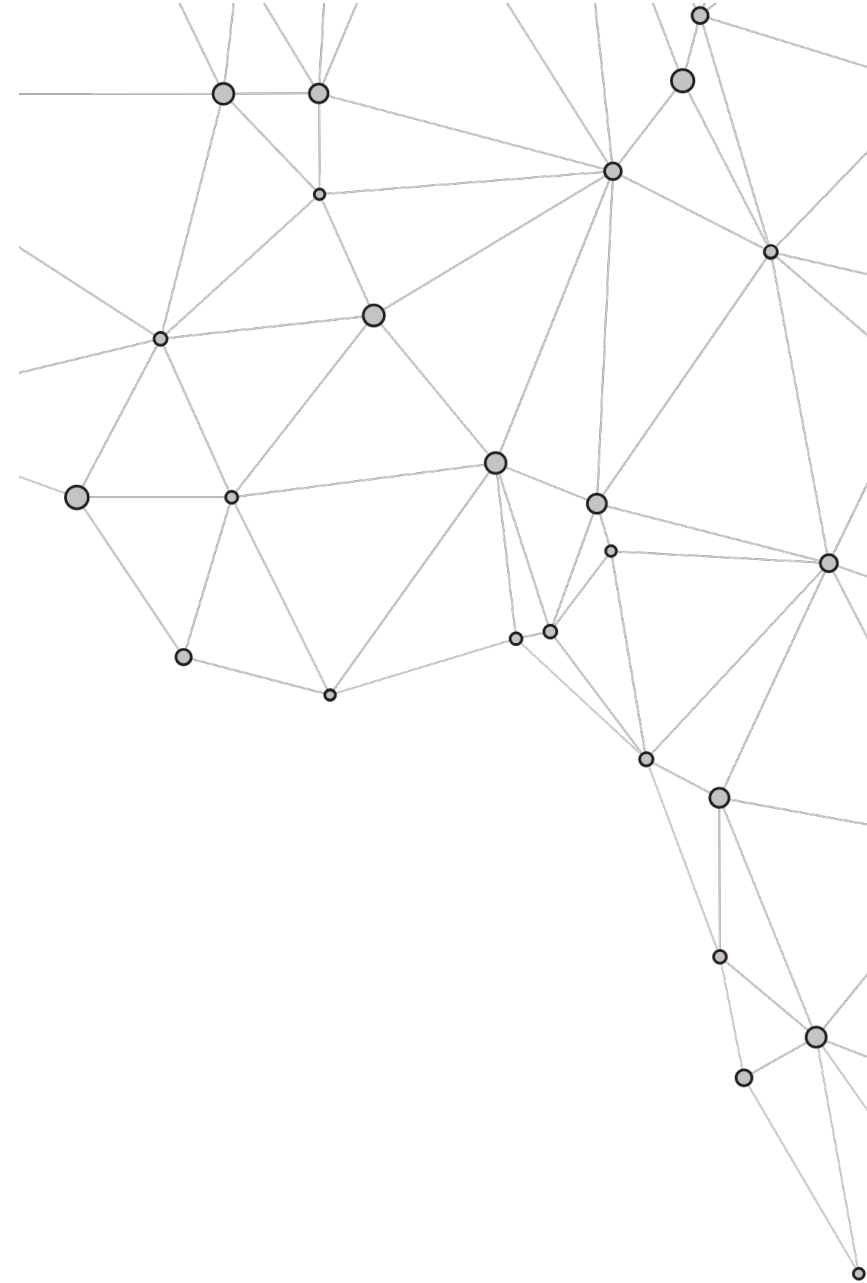
Kết quả

(0, 848) 0.09658593272070018
(0, 1344) 0.1032382690242076
(0, 995) 0.10342093889759825
(0, 1) 0.11196975782384609
(0, 927) 0.11964491432097805
(0, 1562) 0.13307865987371503
(0, 475) 0.10109606028606306
(0, 14) 0.09123904587649263
(0, 967) 0.108680967980731
(0, 883) 0.10506730722446704



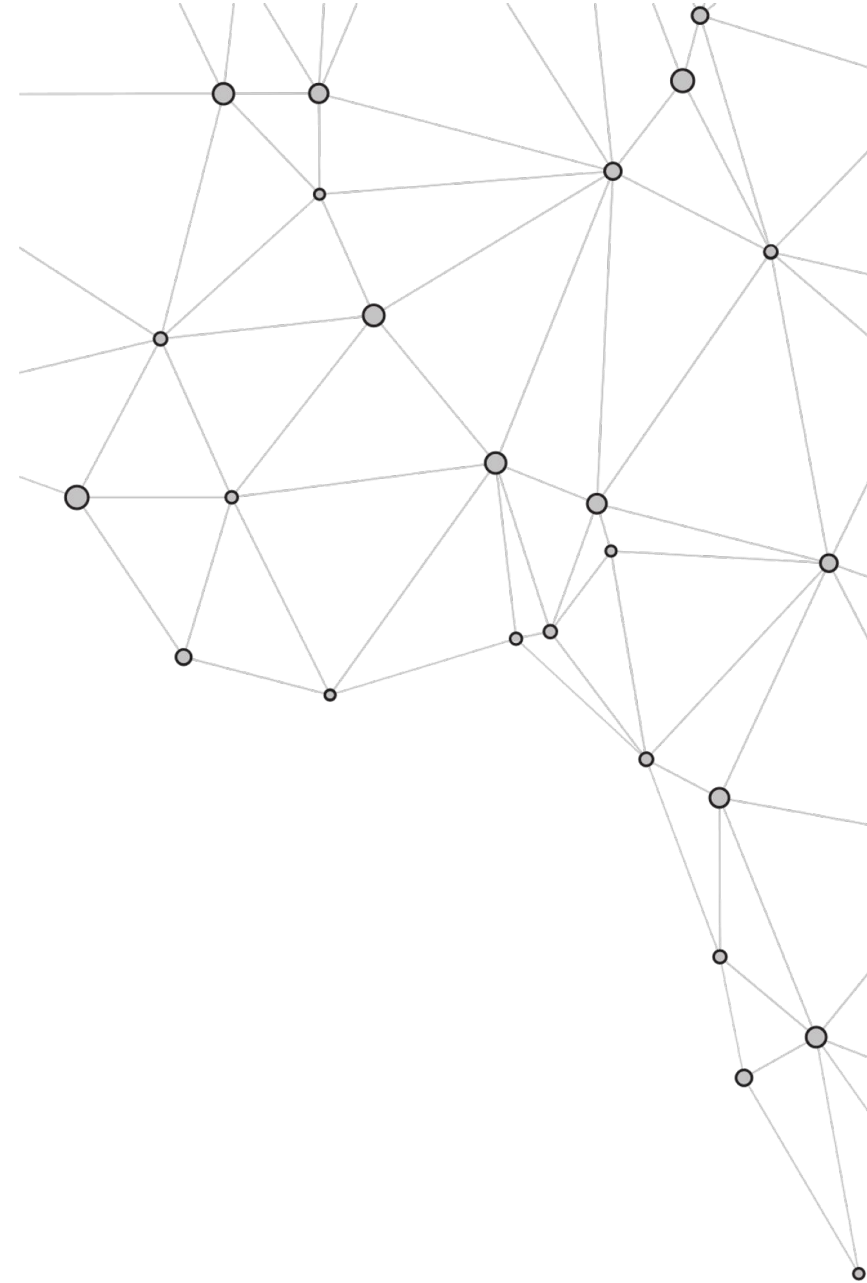
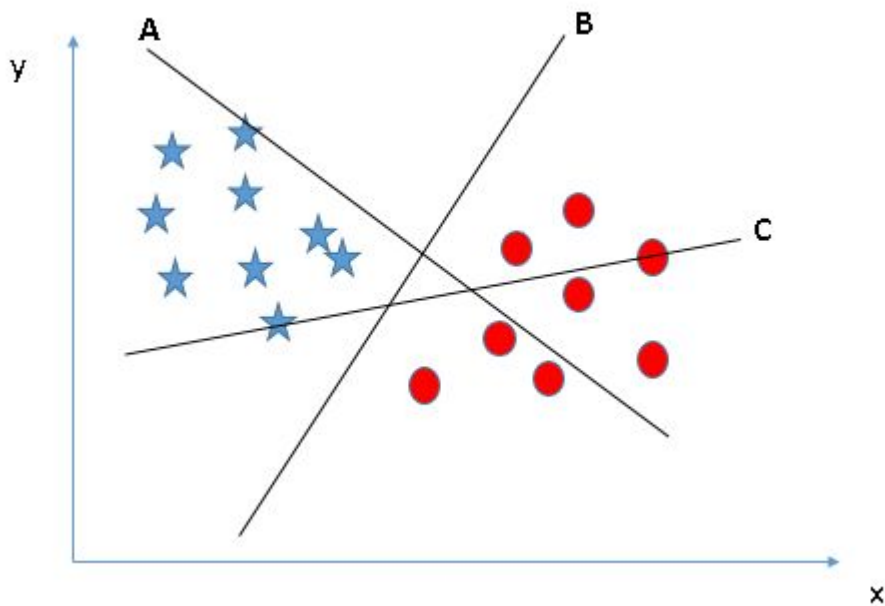
Thuật toán SVM

- **Mục đích:** là thuật toán học có giám sát (**supervised learning**) được sử dụng để phân lớp dữ liệu
- **Ý tưởng:** biểu diễn tập training trong không gian vector, mỗi tài liệu là một điểm. Phương pháp sẽ tìm ra mặt siêu phẳng (**hyperplane**) có thể chia không gian này thành hai lớp riêng biệt.



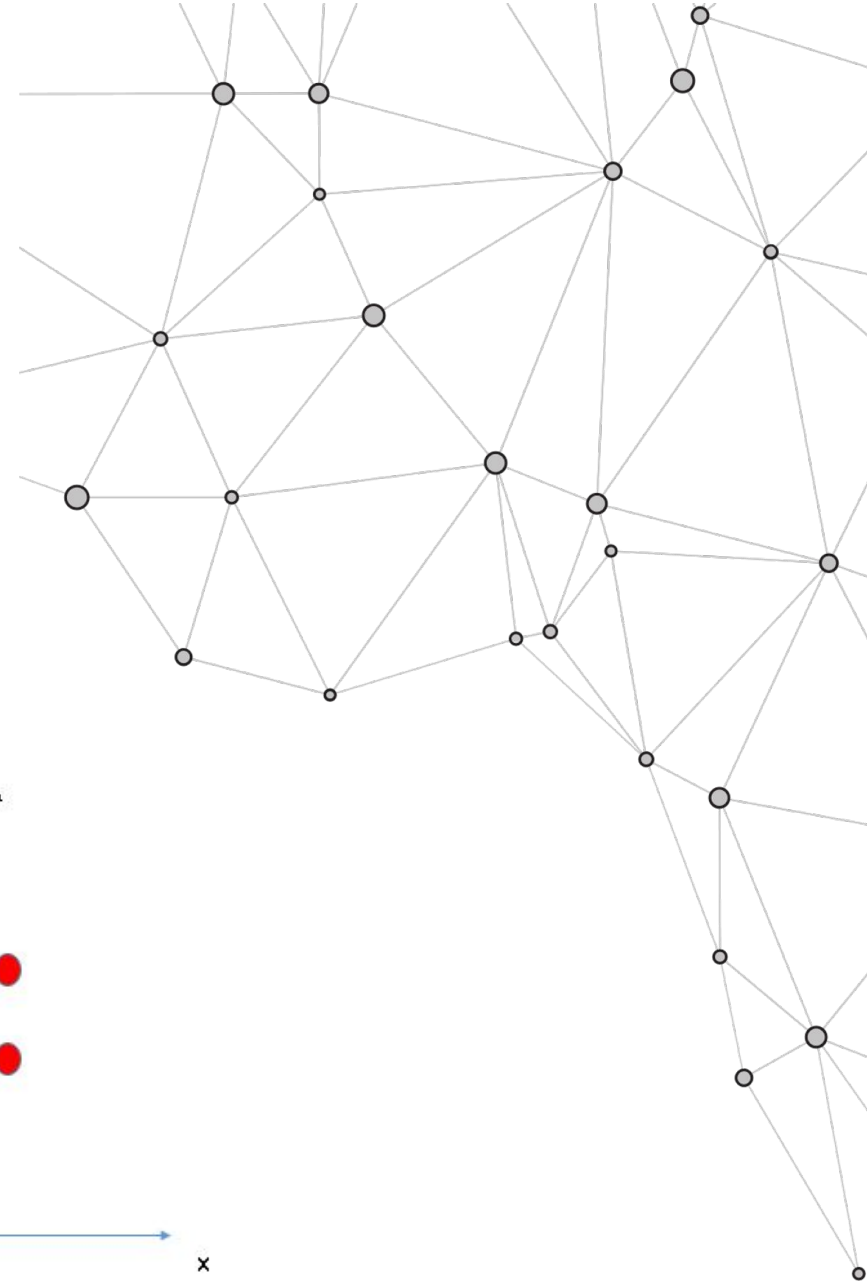
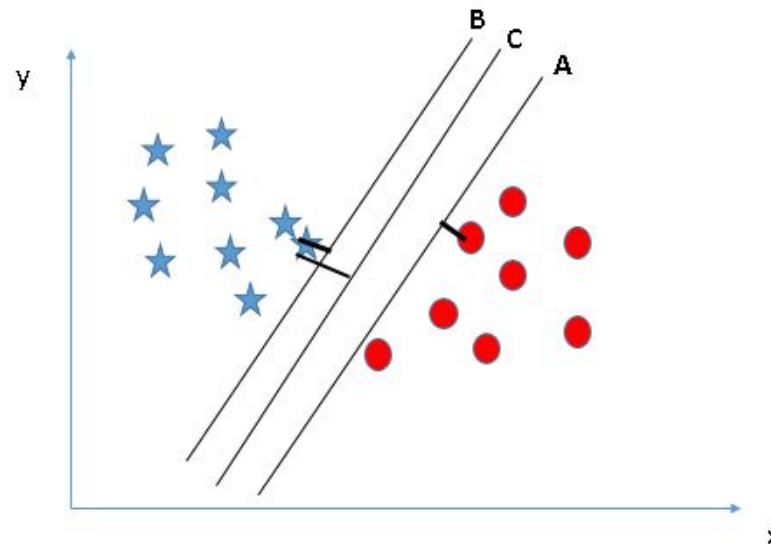
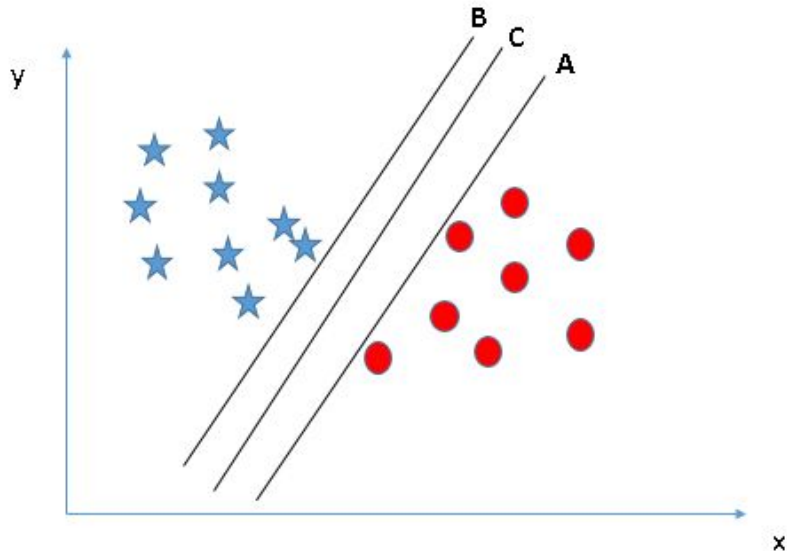
Thuật toán SVM

Quy tắc 1: Chọn 1 siêu phẳng để phân chia thành 2 lớp tốt nhất



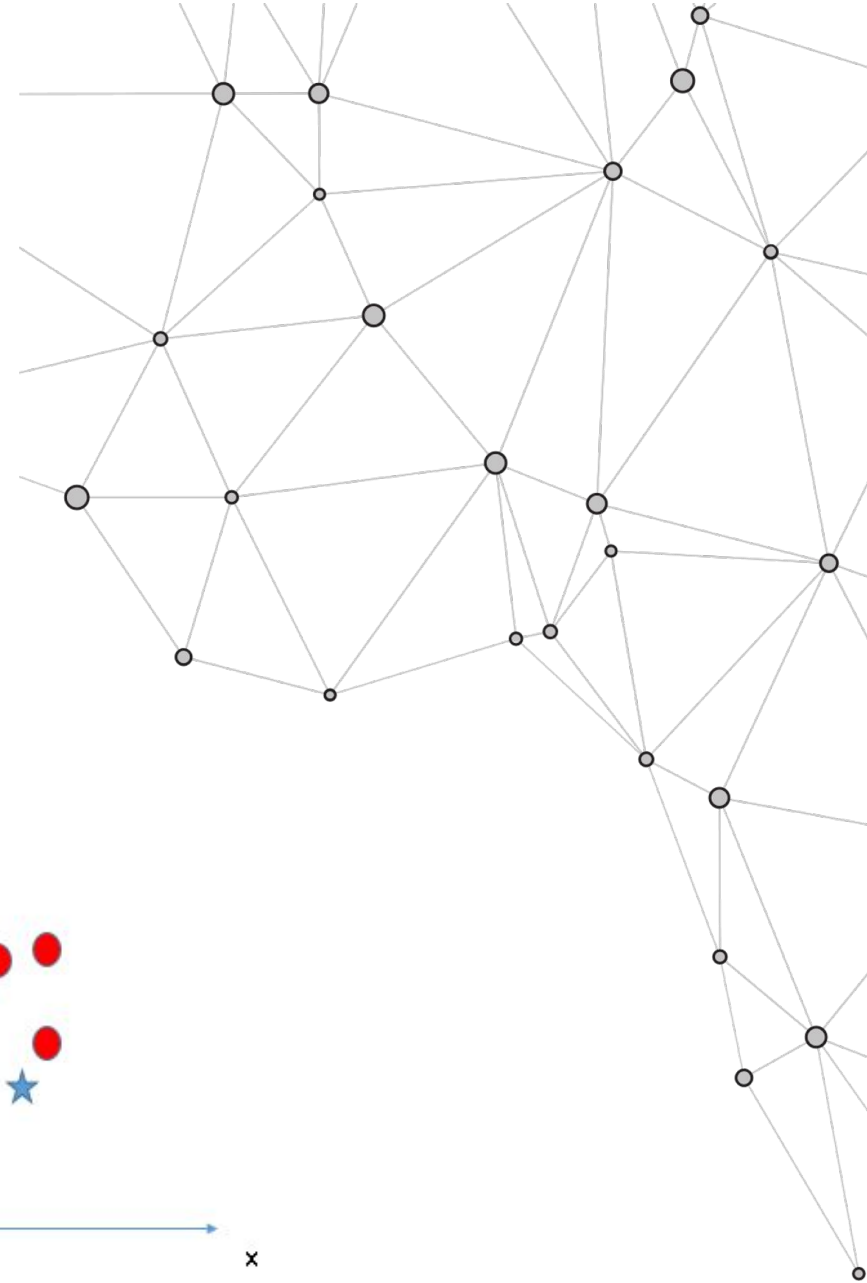
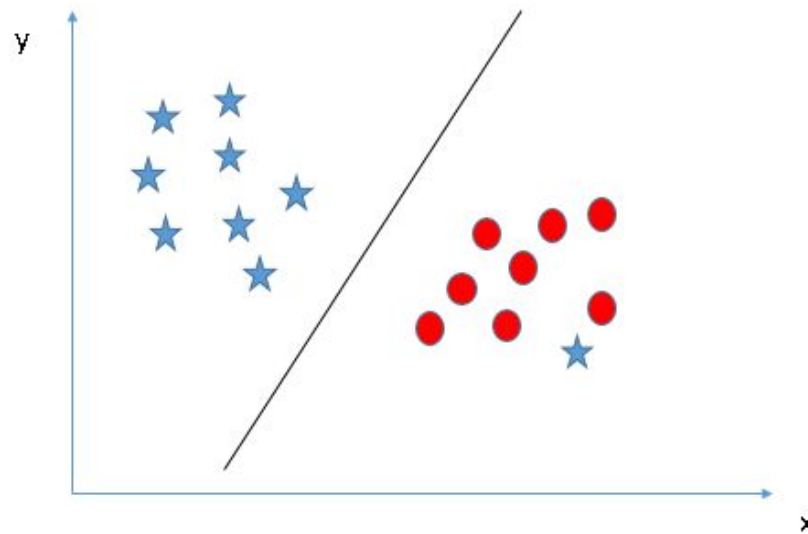
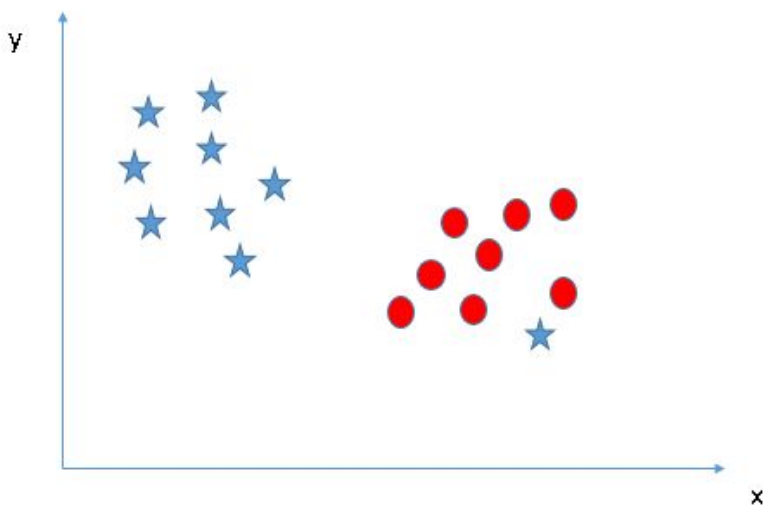
Thuật toán SVM

Quy tắc 2: Xác định khoảng cách lớn nhất giữa điểm dữ liệu gần nhất của hai lớp so với siêu phẳng (Margin)



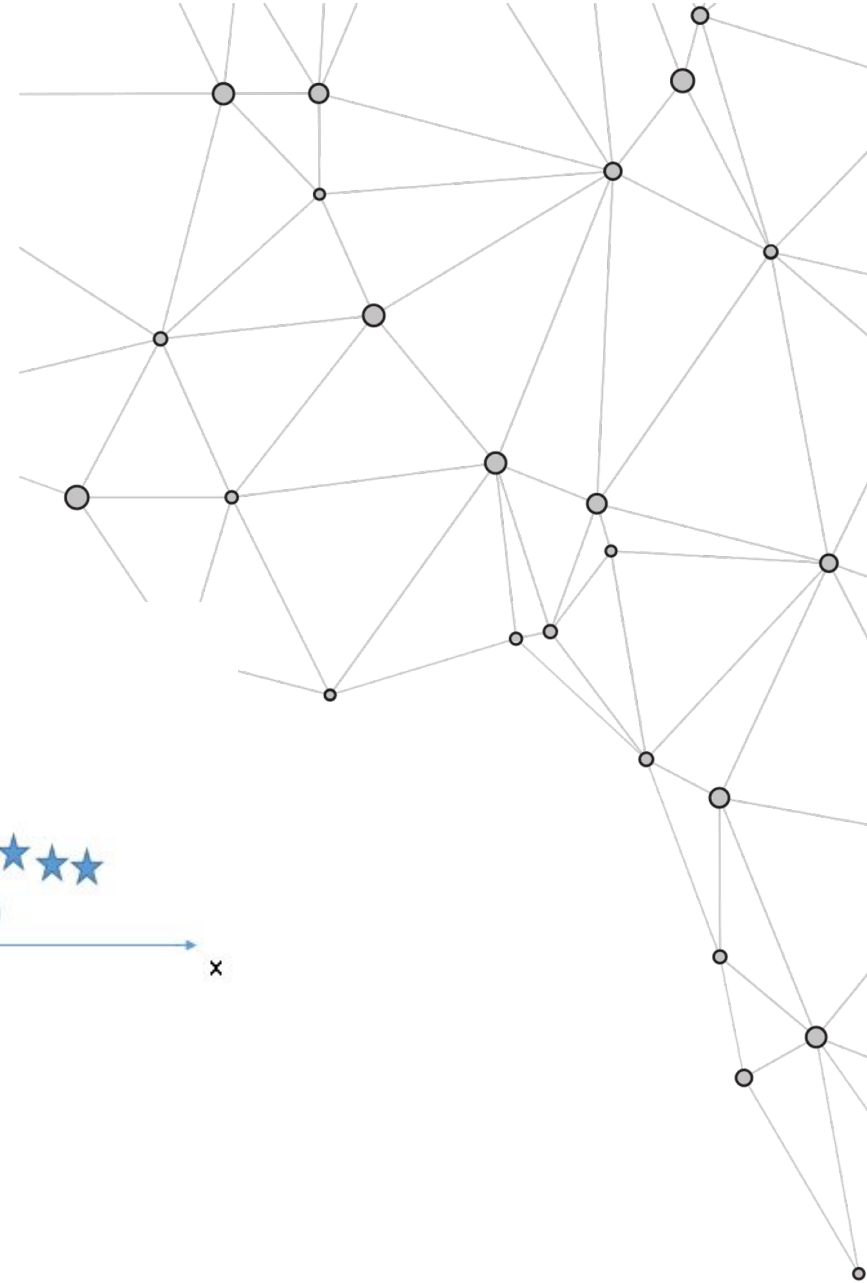
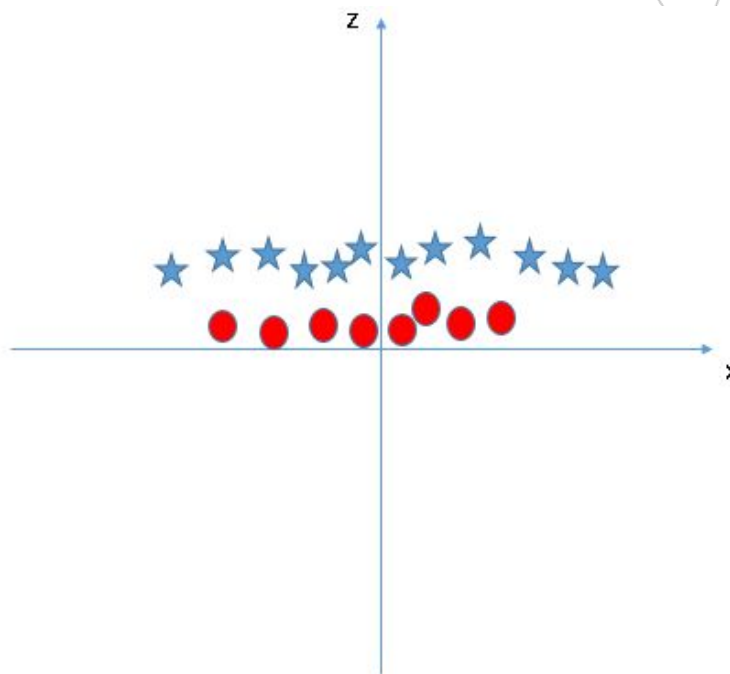
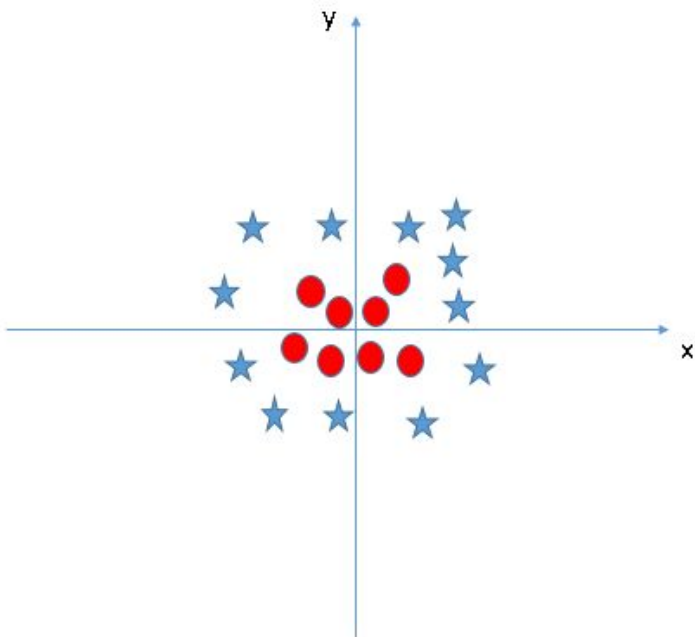
Thuật toán SVM

Quy tắc 3: Chấp nhận ngoại lệ cao



Thuật toán SVM

Quy tắc 4: Thêm tính năng cho SVM



Kết quả thử nghiệm và kết luận

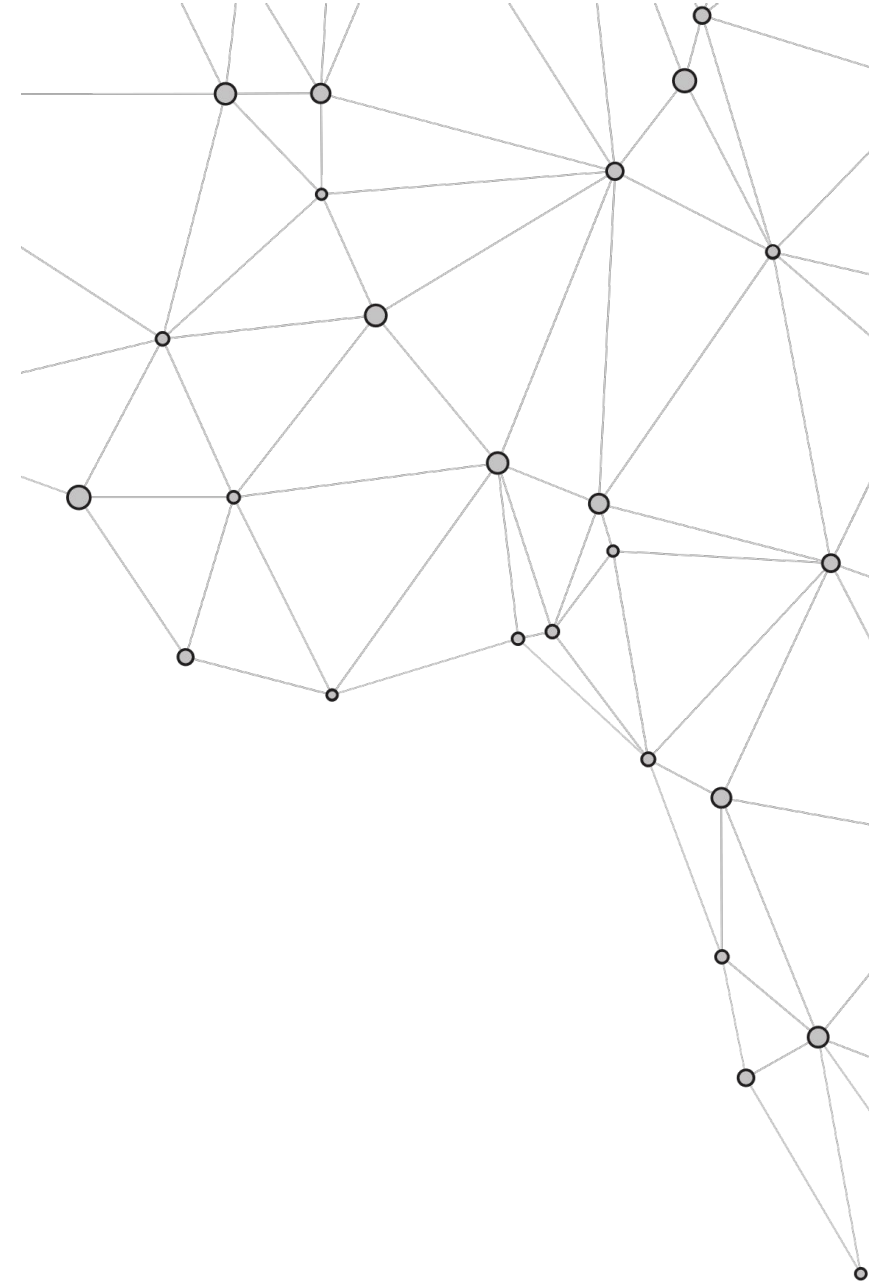
85%

Độ chính xác trung bình

Nhận xét **kết quả**

Nhóm sử dụng hai phương pháp

- Classification Report: giúp hình dung độ chính xác và độ recall tương ứng với từng label, cũng như kết quả tổng thể với toàn bộ label.
- Confusion Matrix: giúp hình dung được lỗi sai thường nằm ở các label nào để có thể tìm được hướng giải quyết.



Thống kê



Kết quả tốt

- __label_Giao_thong (100%)
- __label_Giai_tri (100%)
- __label_May_tinh_va_thiet_bi_dien_tu (100%)



Kết quả thấp

- __label_Tai_chinh (34%)
- __label_Kinh_doanh_va_cong_nghiep (59%)
- __label_Cong_nghe_moi (75%)

Classification report

Measures

	Measures		
	precision	recall	f1-score
weighted avg	0.8574374164863913	0.8714678669667417	0.8569375417086894
macro avg	0.9043794697430698	0.889209689475849	0.892866278621652
accuracy	0.8714678669667417	0.8714678669667417	0.8714678669667417
__label_Thoi_quen_va_so_thich	0.975	0.975	0.975
__label_The_thao	0.9411764705882353	0.8421052631578947	0.8888888888888888
__label_Tai_chinh	0.33884297520661155	0.11884057971014493	0.1759656652360515
__label_Suc_khoe_va_benh_tat	0.8333333333333334	0.8333333333333334	0.8333333333333334
__label_Sach	0.9801980198019802	0.9611650485436893	0.9705882352941178
__label_Phap_luat	0.84	0.6363636363636364	0.7241379310344828
__label_Nha_va_vuon	0.9393939393939394	0.9393939393939394	0.9393939393939394
__label_Nha_dat	0.9751937984496124	0.9889937106918238	0.9820452771272443
__label_Nghe_thuat	0.987012987012987	1.0	0.9934640522875817
__label_Mua_sam	0.9680851063829787	0.934931506849315	0.9512195121951219
__label_May_tinh_va_thiet_bi_dien_tu	1.0	0.94	0.9690721649484536
__label_Mang_internet_va_vien_thong	0.9798657718120806	0.9864864864864865	0.9831649831649831
__label_Lam_dep_va_the_hinh	0.9565217391304348	0.9295774647887324	0.9428571428571428
__label_Kinh_doanh_va_Cong_nghiep	0.5929095354523227	0.8234295415959253	0.689410092395167
__label_Khoa_hoc	0.9512195121951219	0.975	0.9629629629629629
__label_Giao_thong	1.0	1.0	1.0
__label_Giao_duc	0.9378531073446328	0.9540229885057471	0.9458689458689458
__label_Giai_tri	1.0	1.0	1.0
__label_Du_lich	0.9726027397260274	0.9681818181818181	0.970387243735763
__label_Do_an_va_do_uong	0.9782244556113903	0.9915110356536503	0.9848229342327149
__label_Cong_nghe_moi	0.75	0.75	0.75
__label_Con_nguoi_va_xa_hoi	0.9489795918367347	0.9587628865979382	0.9538461538461539
__label_Chinh_tri	0.9543147208121827	0.9447236180904522	0.9494949494949495

Giải thích



Kết quả tốt

Ba nhãn có kết quả cao nhất đều có các cụm từ chuyên ngành đặc biệt và độc nhất với nhãn của mình.

		Measures		
		precision	recall	f1-score
Class	weighted avg	0.8574374164863913	0.8714678669667417	0.8569375417086894
	macro avg	0.9043794697430698	0.889209689475849	0.892866278621652
	accuracy	0.8714678669667417	0.8714678669667417	0.8714678669667417
	__label__Thoi_quen_va_so_thich	0.975	0.975	0.975
	__label__The_thao	0.9411764705882353	0.8421052631578947	0.8888888888888888
	__label__Tai_chinh	0.33884297520661155	0.11884057971014493	0.1759656652360515
	__label__Suc_khoe_va_benh_tat	0.8333333333333334	0.8333333333333334	0.8333333333333334
	__label__Sach	0.9801980198019802	0.9611650485436893	0.9705882352941178
	__label__Phap_luat	0.84	0.6363636363636364	0.7241379310344828
	__label__Nha_va_vuon	0.9393939393939394	0.9393939393939394	0.9393939393939394
	__label__Nha_dat	0.9751937984496124	0.9889937106918238	0.9820452771272443
	__label__Nghe_thuat	0.987012987012987	1.0	0.9934640522875817
	__label__Mua_sam	0.9680851063829787	0.934931506849315	0.9512195121951219
	__label__May_tinh_va_thiet_bi_dien_tu	1.0	0.94	0.9690721649484536
	__label__Mang_internet_va_vien_thong	0.9798657718120806	0.9864864864864865	0.9831649831649831
	__label__Lam_dep_va_the_hinh	0.9565217391304348	0.9295774647887324	0.9428571428571428
	__label__Kinh_doanh_va_Cong_nghiep	0.5929095354523227	0.8234295415959253	0.689410092395167
	__label__Khoa_hoc	0.9512195121951219	0.975	0.9629629629629629
	__label__Giao_thong	1.0	1.0	1.0
	__label__Giao_duc	0.9378531073446328	0.9540229885057471	0.9458689458689458
	__label__Giai_tri	1.0	1.0	1.0
	__label__Du_lich	0.9726027397260274	0.9681818181818181	0.970387243735763
	__label__Do_an_va_do_uong	0.9782244556113903	0.9915110356536503	0.9848229342327149
	__label__Cong_nghe_moi	0.75	0.75	0.75
	__label__Con_nguoi_va_xa_hoi	0.9489795918367347	0.9587628865979382	0.9538461538461539
	__label__Chinh_tri	0.9543147208121827	0.9447236180904522	0.9494949494949495

Giải thích

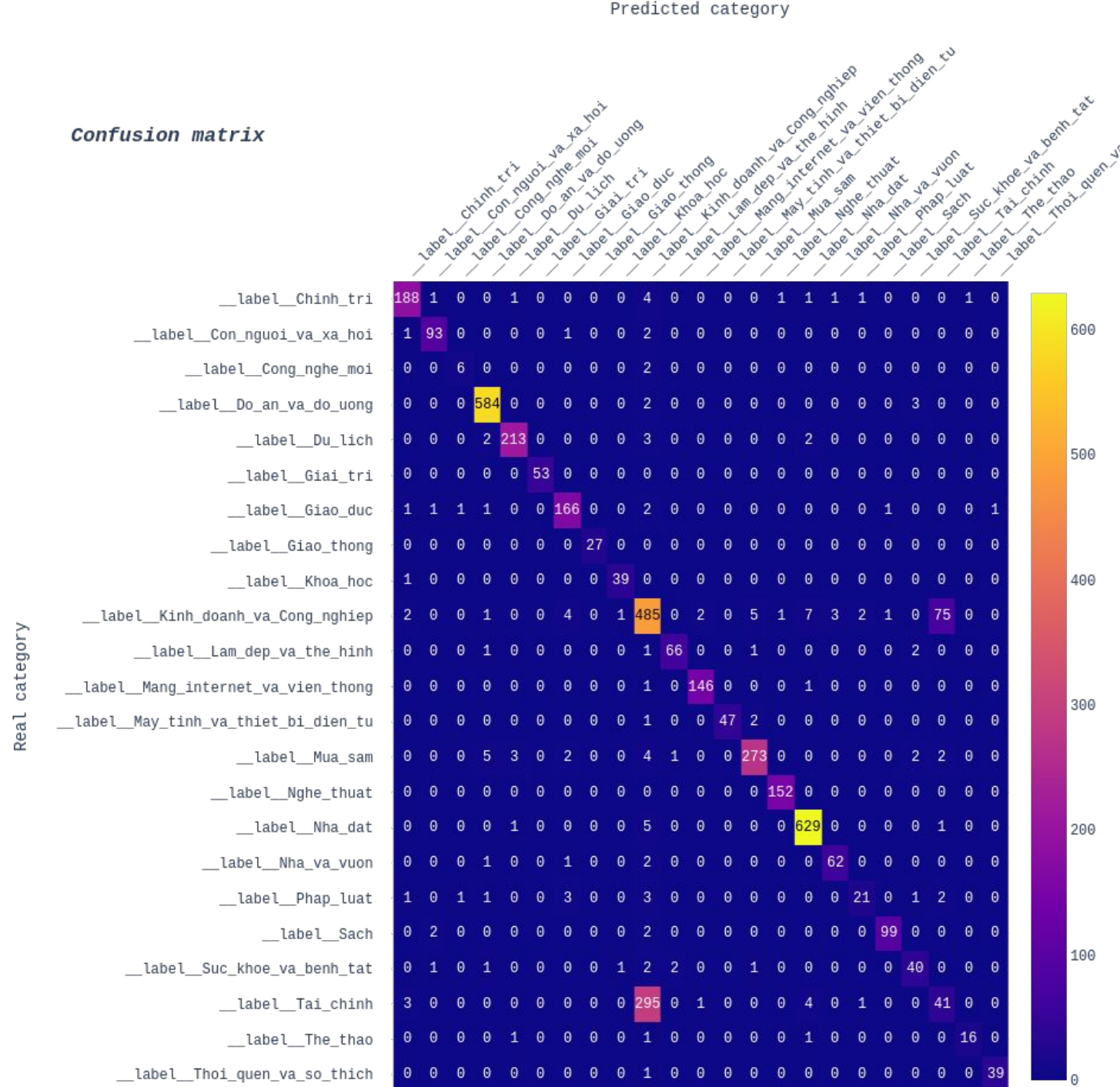
✗ Kết quả thấp

- Với nhãn cong_nghe, kết quả này có thể được giải thích là do nhãn chiếm tỉ lệ rất thấp trong dữ liệu huấn luyện (33 hàng, chiếm 0.2%), nên mô hình không có đủ dữ liệu để đưa ra dự đoán chính xác.
- Với hai nhãn còn lại, độ chính xác thấp là do mô hình thường xuyên dự đoán nhầm tai_chinh thành kinh_doanh và ngược lại. Danh sách từ khóa của hai nhãn có độ tương đồng rất cao, dẫn đến sai sót trong quá trình dự đoán.

Classification report		Measures			
		precision	recall	f1-score	
	weighted avg	0.8574374164863913	0.8714678669667417	0.8569375417086894	
	macro avg	0.9043794697430698	0.889209689475849	0.892866278621652	
	accuracy	0.8714678669667417	0.8714678669667417	0.8714678669667417	
class	__label__Thoi_quen_va_so_thich	0.975	0.975	0.975	
	__label__The_thao	0.9411764705882353	0.8421052631578947	0.8888888888888888	
	__label__Tai_chinh	0.33884297520661155	0.11884057971014493	0.1759656652360515	
	__label__Suc_khoe_va_benh_tat	0.8333333333333334	0.8333333333333334	0.8333333333333334	
	__label__Sach	0.9801980198019802	0.9611650485436893	0.9705882352941178	
	__label__Phap_luat	0.84	0.6363636363636364	0.7241379310344828	
	__label__Nha_va_vuon	0.9393939393939394	0.9393939393939394	0.9393939393939394	
	__label__Nha_dat	0.9751937984496124	0.9889937106918238	0.9820452771272443	
	__label__Nghe_thuat	0.987012987012987	1.0	0.9934640522875817	
	__label__Mua_sam	0.9680851063829787	0.934931506849315	0.9512195121951219	
	__label__May_tinh_va_thiet_bi_dien_tu	1.0	0.94	0.9690721649484536	
	__label__Mang_internet_va_vien_thong	0.9798657718120806	0.9864864864864865	0.9831649831649831	
	__label__Lam_dep_va_the_hinh	0.9565217391304348	0.9295774647887324	0.9428571428571428	
	__label__Kinh_doanh_va_Cong_nghiep	0.5929095354523227	0.8234295415959253	0.689410092395167	
	__label__Khoa_hoc	0.9512195121951219	0.975	0.9629629629629629	
	__label__Giao_thong	1.0	1.0	1.0	
	__label__Giao_duc	0.9378531073446328	0.9540229885057471	0.9458689458689458	
	__label__Giai_tri	1.0	1.0	1.0	
	__label__Du_lich	0.9726027397260274	0.9681818181818181	0.970387243735763	
	__label__Do_an_va_do_uong	0.9782244556113903	0.9915110356536503	0.9848229342327149	
	__label__Cong_nghe_moi	0.75	0.75	0.75	
	__label__Con_nguoi_va_xa_hoi	0.9489795918367347	0.9587628865979382	0.9538461538461539	
	__label__Chinh_tri	0.9543147208121827	0.9447236180904522	0.9494949494949495	

Kết quả thấp

Kết quả thấp



Kết luận

Tóm lại, các kết quả nói trên đã phản ánh được bản chất của các chủ đề cũng như hạn chế của phương pháp sử dụng. Các chủ đề thường không độc lập với nhau mà cũng bao hàm và giao nhau, dẫn tới sai sót trong quá trình dự đoán. Mặt khác, phương pháp sử dụng không quan tâm tới cấu trúc câu mà chỉ tập trung vào số lần xuất hiện của từ cũng một phần làm giảm độ chính xác của mô hình.

