

TỔNG LIÊN ĐOÀN LAO ĐỘNG VIỆT NAM  
TRƯỜNG ĐẠI HỌC TÔN ĐỨC THẮNG  
KHOA CÔNG NGHỆ THÔNG TIN



ĐỒ ÁN  
HỆ THỐNG THƯƠNG MẠI THÔNG MINH

# XÂY DỰNG MÔ HÌNH DỰ ĐOÁN TẬP KHÁCH HÀNG NGỪNG SỬ DỤNG DỊCH VỤ CỦA NGÂN HÀNG

*Người hướng dẫn:* Th.S DƯƠNG HỮU PHÚC

*Người thực hiện:* LÊ NGUYỄN MINH TUẤN - 51800950

TỪ HUY VẠN - 51800263

*Lớp:* 18050301

*Khóa:* 22

THÀNH PHỐ HỒ CHÍ MINH, NĂM 2021

TỔNG LIÊN ĐOÀN LAO ĐỘNG VIỆT NAM  
TRƯỜNG ĐẠI HỌC TÔN ĐỨC THẮNG  
KHOA CÔNG NGHỆ THÔNG TIN



ĐỒ ÁN  
HỆ THỐNG THƯƠNG MẠI THÔNG MINH

# XÂY DỰNG MÔ HÌNH DỰ ĐOÁN TẬP KHÁCH HÀNG NGỪNG SỬ DỤNG DỊCH VỤ CỦA NGÂN HÀNG

*Người hướng dẫn:* Th.S DƯƠNG HỮU PHÚC

*Người thực hiện:* LÊ NGUYỄN MINH TUẤN - 51800950

TỪ HUY VẠN - 51800263

*Lớp:* 18050301

*Khóa:* 22

THÀNH PHỐ HỒ CHÍ MINH, NĂM 2021

## LỜI CẢM ƠN

Qua một học kỳ đầy khó khăn và vất vả, chúng em đã học tập được nhiều thứ từ môn 'Hệ thống thương mại thông minh'. Với sự giúp đỡ chỉ bảo và giảng dạy nhiệt tình của quý thầy cô, đặc biệt là quý cô Khoa Công nghệ thông tin Trường Đại học Tôn Đức Thắng.

Sau quãng thời gian tìm hiểu và thực hiện, chúng em đã hoàn thành đồ án cuối kỳ môn 'Hệ thống thương mại thông minh'. Em xin dành lời cảm ơn đặc biệt đến Thầy Dương Hữu Phúc đã giúp đỡ và chỉ bảo những thiếu sót trong quá trình thực hiện của chúng em.

Và đặc biệt, trong học kỳ này, Khoa đã tổ chức cho chúng em được tiếp cận với môn học mà theo hai em là rất hữu ích đối với sinh viên ngành Khoa Học Máy Tính cũng như tất cả các sinh viên thuộc các chuyên ngành Công Nghệ Thông Tin khác. Chúng em xin chân thành cảm ơn thầy Dương Hữu Phúc đã tận tâm hướng dẫn chúng em qua từng buổi học trên lớp thảo luận về lĩnh vực sáng tạo trong nghiên cứu khoa học.

Trong quá trình làm tiểu luận, khó tránh khỏi những sai sót, rất mong Thầy bỏ qua. Đồng thời em rất mong nhận được những ý kiến đóng góp của Thầy để chúng em có thể tiếp thu và cải thiện để nâng cao kiến thức của mình, đồng thời phục vụ tốt hơn cho quá trình nghiên cứu sau này.

Chúng em xin chân thành cảm ơn!

## ĐỒ ÁN ĐƯỢC HOÀN THÀNH TẠI TRƯỜNG ĐẠI HỌC TÔN ĐỨC THẮNG

Chúng tôi xin cam đoan đây là thành quả của riêng chúng tôi và được sự hướng dẫn của ThS Dương Hữu Phúc. Các nội dung nghiên cứu, kết quả trong đề tài này là trung thực và chưa công bố dưới bất kỳ hình thức nào trước đây. Những số liệu trong các bảng biểu phục vụ cho việc phân tích, nhận xét, đánh giá được chính tác giả thu thập từ các nguồn khác nhau có ghi rõ trong phần tài liệu tham khảo.

Ngoài ra, trong đồ án còn sử dụng một số nhận xét, đánh giá cũng như số liệu của các tác giả khác, cơ quan tổ chức khác đều có trích dẫn và chú thích nguồn gốc.

**Nếu phát hiện có bất kỳ sự gian lận nào chúng tôi xin hoàn toàn chịu trách nhiệm về nội dung đồ án của mình.** Trường đại học Tôn Đức Thắng không liên quan đến những vi phạm tác quyền, bản quyền do chúng tôi gây ra trong quá trình thực hiện (nếu có).

TP. Hồ Chí Minh, ngày            tháng            năm

*Tác giả*

*(ký tên và ghi rõ họ tên)*

*Lê Nguyễn Minh Tuấn*

*Từ Huy Vạn*

# PHẦN NHẬN XÉT VÀ ĐÁNH GIÁ CỦA GIẢNG VIÊN

Phần xác nhận của GV hướng dẫn

---

---

---

---

---

---

---

---

TP. Hồ Chí Minh, ngày            tháng            năm  
(ký tên và ghi rõ họ tên)

Phần đánh giá của GV chấm bài

---

---

---

---

---

---

---

---

TP. Hồ Chí Minh, ngày            tháng            năm  
(ký tên và ghi rõ họ tên)

## LỜI MỞ ĐẦU

Với sự phát triển rất nhanh của xã hội hiện nay, định nghĩa về cách mạng khoa học kỹ thuật 4.0 không phải là điều quá xa lạ. Với lĩnh vực trí tuệ nhân tạo là một trong những trọng tâm chính, nó đang ảnh hưởng đến tất cả các lĩnh vực, nền kinh tế, các ngành kinh tế và ngành công nghiệp. Ngân hàng là một lĩnh vực lâu đời, bắt nguồn từ thế kỷ 13 và tại đây trong thế kỷ 21 trí tuệ nhân tạo đang bắt đầu đặt những ảnh hưởng đầu tiên của mình lên lĩnh vực lâu đời này.

Ngân hàng giờ đây là một trong những lĩnh vực cạnh tranh, sự cạnh tranh gay gắt này buộc các ngân hàng phải cải thiện chất lượng của mình nếu muốn giữ lấy khách hàng và trí tuệ nhân tạo với các mô hình học máy là một trong những công cụ hiệu quả nhất để làm được điều đó.

Trí tuệ nhân tạo hay còn được gọi là Artificial Intelligence được hiểu như một ngành của khoa học máy tính liên quan đến việc tự động hóa các hành vi thông minh. AI là trí tuệ do con người lập trình tạo nên với mục tiêu giúp máy tính có thể tự động hóa các hành vi thông minh như con người. Trí tuệ nhân tạo khác với việc lập trình logic trong các ngôn ngữ lập trình là ở việc ứng dụng các hệ thống học máy để mô phỏng trí tuệ con người trong các xử lý mà con người làm tốt hơn máy tính. Cụ thể, trí tuệ nhân tạo giúp máy tính có được những trí tuệ của con người như biết suy nghĩ và lập luận để giải quyết vấn đề, biết giao tiếp do hiểu ngôn ngữ, tiếng nói, biết học và tự thích nghi, v.v...

Với sức mạnh tính toán vượt xa con người cùng với trí thông minh, trí tuệ nhân tạo là một công cụ tuyệt vời để thống kê và thực hiện các phép toán phức tạp từ đó giúp con người có nhiều nguồn thông tin hơn để vạch ra những chiến lược cụ thể trong khối ngành kinh tế. Và lĩnh vực ngân hàng không phải ngoại lệ.

Trong phạm vi đề án, chúng em sẽ trực quan hóa các yếu tố ảnh hưởng đến quyết định của khách hàng khi ra quyết định liệu có ngưng sử dụng dịch vụ của ngân hàng hay không và từ các yếu tố đó xây dựng các mô hình Machine Learning để dự đoán xem với các yếu tố từ thông tin của khách hàng, liệu khách hàng có ngưng

sử dụng dịch vụ hay vẫn tiếp tục gắn bó với ngân hàng.

Việc làm đồ án là giúp cho chúng em có cơ hội áp dụng, và được tổng hợp các kiến thức của chúng em đã được học trên lớp. Nhưng với vốn kinh nghiệm ít ỏi và kiến thức chưa sâu, sai sót là không thể tránh khỏi. Rất mong có thể nhận được những đánh giá và chỉ bảo nhiệt tình của quý Thầy, Cô.

Chúng em xin cảm ơn thầy Dương Hữu Phúc đã nhiệt tình hướng dẫn và hoàn thiện đồ án đúng hạn !

# Mục lục

<b>1</b>	<b>TỔNG QUAN ĐỀ TÀI</b>	<b>1</b>
1.1	GIỚI THIỆU ĐỀ TÀI . . . . .	1
1.2	PHÁT BIỂU BÀI TOÁN . . . . .	1
1.3	MỤC TIÊU ĐỀ TÀI . . . . .	2
1.4	PHẠM VI ĐỀ TÀI . . . . .	2
1.5	Ý NGHĨA KHOA HỌC & THỰC TIỄN . . . . .	2
1.5.1	Ý nghĩa khoa học . . . . .	2
1.5.2	Ý nghĩa thực tiễn . . . . .	3
1.6	CẤU TRÚC BÁO CÁO . . . . .	3
<b>2</b>	<b>TỔNG QUAN VỀ GIẢI THUẬT</b>	<b>4</b>
2.1	LOGISTICS REGRESSION . . . . .	4
2.1.1	Tổng quan về Logisctics Regression . . . . .	4
2.1.2	Model Implement . . . . .	5
2.2	K-NEAREST NEIGBORS . . . . .	5
2.2.1	Tổng quan về K-Nearest Neighbors . . . . .	5
2.2.2	Model Implement . . . . .	7
2.3	ADAPTIVE BOOSTING . . . . .	7



2.3.1	Tổng quan về Adaptive Boosting . . . . .	7
2.3.2	Model Implement . . . . .	8
2.4	RANDOM FOREST . . . . .	9
2.4.1	Tổng quan về Random Forest . . . . .	9
2.4.2	Model Implement . . . . .	10
2.5	GRADIENT BOOSTING . . . . .	10
2.5.1	Tổng quan về Gradient Boosting . . . . .	10
2.5.2	Model Implement . . . . .	12
<b>3</b>	<b>DỮ LIỆU THỰC NGHIỆM</b>	<b>13</b>
3.1	PHÂN TÍCH BỘ DỮ LIỆU . . . . .	13
3.1.1	Nguồn gốc data . . . . .	13
3.1.2	Chi tiết các thuộc tính . . . . .	13
3.1.3	Trình bày bộ dữ liệu bằng Tableau . . . . .	14
3.2	TIỀN XỬ LÝ DỮ LIỆU . . . . .	20
3.3	CÁC PHƯƠNG PHÁP ĐÁNH GIÁ GIẢI THUẬT . . . . .	22
3.3.1	Giới thiệu về các phương pháp đánh giá . . . . .	22
3.3.2	Accuracy Score . . . . .	22
3.3.3	Confusion Matrix . . . . .	23
3.3.4	Precision và Recall . . . . .	24
3.3.5	F1-Score . . . . .	25
3.4	TỔNG QUAN VỀ SMOTE . . . . .	25
3.4.1	SMOTE là gì ? . . . . .	25
3.4.2	Tại sao lại cần sử dụng SMOTE ? . . . . .	25

<b>4</b>	<b>THỰC NGHIỆM</b>	<b>27</b>
4.1	IMPORT CÁC THƯ VIỆN VÀ ĐỌC DỮ LIỆU . . . . .	27
4.2	ĐÁNH GIÁ TRÊN CÁC MÔ HÌNH . . . . .	28
4.2.1	Logistic Regression . . . . .	28
4.2.2	K-Nearest Neighbors . . . . .	32
4.2.3	Adaptive Boosting . . . . .	36
4.2.4	Random Forest . . . . .	40
4.2.5	Gradient Boosting . . . . .	44
<b>5</b>	<b>KẾT LUẬN</b>	<b>49</b>

# Danh sách hình vẽ

2.1	Biểu diễn Logistic Regression theo Neural Network . . . . .	4
2.2	Các cách cơ bản để tính khoảng cách giữa 2 điểm dữ liệu . . . . .	6
2.3	Các bước hoạt động của Adaptive Boosting . . . . .	8
2.4	Các bước hoạt động của Random Forest . . . . .	9
2.5	Cách hoạt động của Gradient Boosting . . . . .	11
3.1	Churn vs Not Churn . . . . .	15
3.2	Gender . . . . .	16
3.3	Age group . . . . .	17
3.4	Income . . . . .	18
3.5	Churn Data . . . . .	19
3.6	Data Overview . . . . .	20
3.7	10 cột đầu tiên trước khi tiền xử lý . . . . .	21
3.8	10 cột đầu tiên sau khi tiền xử lý . . . . .	21
4.1	10 cột đầu tiên của dữ liệu giả dùng để đánh giá . . . . .	28
4.2	Đánh giá Logistics Regression qua confusion matrix không sử dụng SMOTE . . . . .	30

4.3	Đánh giá Logistics Regression qua confusion matrix có sử dụng SMOTE . . . . .	32
4.4	Đánh giá KNN qua confusion matrix không sử dụng SMOTE . . . .	34
4.5	Đánh giá KNN qua confusion matrix có sử dụng SMOTE . . . . .	36
4.6	Đánh giá Adaptive Boosting qua confusion matrix không sử dụng SMOTE . . . . .	38
4.7	Đánh giá Adaptive Boosting qua confusion matrix có sử dụng SMOTE	40
4.8	Đánh giá Random Forest qua confusion matrix không sử dụng SMOTE	42
4.9	Đánh giá Random Forest qua confusion matrix có sử dụng SMOTE	44
4.10	Đánh giá Gradient Boosting qua confusion matrix không sử dụng SMOTE . . . . .	46
4.11	Đánh giá Gradient Boosting qua confusion matrix có sử dụng SMOTE	48

## BẢNG PHÂN CÔNG CÔNG VIỆC

MSSV	HỌ TÊN SINH VIÊN	CÔNG VIỆC
51800950	Lê Nguyễn Minh Tuấn	Report & Coding
51800263	Từ Huy Vạn	Report & Tableau

Bảng 1: Bảng phân công công việc

## Chương 1

# TỔNG QUAN ĐỀ TÀI

### 1.1 GIỚI THIỆU ĐỀ TÀI

Với việc chi phí để có được những khách hàng mới ngày càng gia tăng nên vì vậy việc duy trì các khách hàng cũ trở nên vô cùng quan trọng. Vì vậy các ngân hàng cần phải dự đoán trước các khách hàng nào đang có nguy cơ rời bỏ ngân hàng để kịp thời có chiến lược giữ chân khách hàng.

Bộ dữ liệu được nghiên cứu bao gồm 10000 khách hàng đang sử dụng và ngưng sử dụng dịch vụ ngân hàng từ độ tuổi 26 đến 73 bao gồm nam và nữ.

### 1.2 PHÁT BIỂU BÀI TOÁN

Bài toán đặt ra là làm cách nào có thể dự đoán được khách hàng nào có khả năng ngừng sử dụng dịch vụ ngân hàng trong tương lai hay không. Các yếu tố nào làm ảnh hưởng tới quyết định của khách hàng, điều này thôi thúc các doanh nghiệp tìm ra các phương pháp để dự đoán hành vi của khách hàng dựa trên dữ liệu có sẵn và các phương pháp Machine Learning là một công cụ hiệu quả để giải quyết các vấn đề trên.

1. Các mô hình Machine Learning sẽ tính xác suất khách hàng có khả năng bỏ đi dựa theo dữ liệu truyền vào của khách hàng đó.

2. Phân tích được thuộc tính nào của khách hàng giữ vai trò trọng yếu trong quyết định cuối cùng của khách hàng

## 1.3 MỤC TIÊU ĐỀ TÀI

Vậy Machine Learning là gì ?

Machine learning hay học máy là một kỹ thuật cho phép máy tính có khả năng cải thiện chính bản thân chúng dựa trên dữ liệu mẫu (training data) hoặc dựa vào kinh nghiệm (những gì đã được học). Bằng cách áp dụng Machine learning có thể dự đoán được khách nào đang có nguy cơ bỏ đi.

Dựa vào bộ dữ liệu kèm theo, trực quan hóa dữ liệu bằng Tableau từ đó phân tích được cách yếu tố ảnh hưởng tới quyết định của khách hàng kết hợp với việc huấn luyện các mô hình Machine Learning để cuối cùng có những công cụ hiệu quả dự đoán được hành vi của khách hàng

## 1.4 PHẠM VI ĐỀ TÀI

Các thuật toán Machine Learning được áp dụng trong đề tài:

1. Logistics Regression
2. K-Nearest Neighbors
3. Adaptive Boosting
4. Random Forest
5. Gradient Boosting

Các thuật toán này được áp dụng để phân tích trong bộ dữ liệu kèm theo, liên quan đến lĩnh vực ngân hàng và mô hình để dự đoán là phân loại nhị phân. Nghĩa là khách hàng vẫn sẽ sử dụng dịch vụ ngân hàng hoặc ngưng sử dụng

## 1.5 Ý NGHĨA KHOA HỌC & THỰC TIỄN

### 1.5.1 Ý nghĩa khoa học

Đề tài này sẽ giúp làm sáng tỏ đặc điểm và thực trạng của việc khách hàng từ bỏ dịch vụ ngân hàng. Bằng cách nắm bắt và làm rõ các khía

cạnh chính trong mối quan hệ giữa khách hàng và ngân hàng ta có thể xây dựng một kho dữ liệu cục bộ về khách hàng. Điều này sẽ làm nền tảng không chỉ cho việc dự đoán rời đi của khách hàng mà còn giúp trong những vấn đề khác như là tính toán giá trị lâu dài của khách hàng, gợi ý các dịch vụ phù hợp cho khách hàng.

### **1.5.2 Ý nghĩa thực tiễn**

Nếu có thể dự đoán được khách hàng nào chuẩn bị ngừng sử dụng ngân hàng thì ngân hàng có thể thay đổi các chiến lược cũng như các ưu đãi để giữ chân khách hàng. Như vậy ngân hàng có thể phát triển mối quan hệ với khách hàng mới và giữ chân những khách hàng cũ nhằm tiết kiệm chi phí hơn. Đồng thời việc phân tích dữ liệu trên sẽ hỗ trợ cho ngân hàng cải thiện về mức độ phục vụ khách hàng.

## **1.6 CẤU TRÚC BÁO CÁO**

Báo cáo được chia thành các chương như sau :

Chương 1 : Tổng quan đề tài

Chương 2 : Tổng quan về giải thuật

Chương 3 : Dữ liệu thực nghiệm

Chương 4 : Thực nghiệm

Chương 5 : Kết luận

Chương 6 : Tài liệu tham khảo

Trong quá trình làm bài, chúng em đã cố gắng trình bày một cách đầy đủ, dễ hiểu nhất, đồng thời chúng em đã nỗ lực trong quá trình tìm hiểu tài liệu và viết bài tiểu luận này, nhưng không thể tránh khỏi những sai lầm, thiếu sót được. Chúng em, mong sẽ nhận được ý kiến đóng góp từ phía thầy cô.



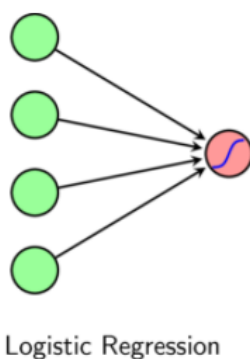
## Chương 2

# TỔNG QUAN VỀ GIẢI THUẬT

### 2.1 LOGISTICS REGRESSION

#### 2.1.1 Tổng quan về Logisctics Regression

Logistic Regression là 1 thuật toán phân loại được dùng để gán các đối tượng cho 1 tập hợp giá trị rời rạc (như 0, 1, 2, ...). Một ví dụ điển hình là phân loại Email, gồm có email công việc, email gia đình, email spam, ... Giao dịch trực tuyến có là an toàn hay không an toàn, khối u lành tính hay ác tính. Thuật toán trên dùng hàm sigmoid logistic để đưa ra đánh giá theo xác suất. Ví dụ: Khối u này 80% là lành tính, giao dịch này 90% là gian lận, ...



Hình 2.1: Biểu diễn Logistic Regression theo Neural Network

**Ưu điểm**

- Đơn giản, hiệu quả và được sử dụng rộng rãi
- Không yêu cầu mở rộng các tính năng
- Cung cấp điểm xác suất để quan sát
- Không đòi hỏi sức mạnh tính toán cao
- Tính toán nhanh trong việc phân loại các biến không xác định

**Nhược điểm**

- Nếu số lượng data ít hơn số lượng thuộc tính thì sẽ xảy ra overfitting
- Không thể giải quyết vấn đề phi tuyến tính
- Không hoạt động tốt với các biến độc lập không tương quan với biến mục tiêu
- Khó tính toán các mối quan hệ phức tạp, các thuật toán như Neural Networks cho kết quả tốt hơn Logistic Regression
- Yêu cầu các điểm dữ liệu được tạo ra một cách độc lập với nhau

**2.1.2 Model Implement**

```

1 # import Logistics Regression
2 from sklearn.linear_model import LogisticRegression
3
4 # model implement
5 lr = LogisticRegression(random_state=42, solver='liblinear', multi_class='ovr')
6 lr.fit(X_train,y_train)
7
8 # predict from model
9 y_pred = lr.predict(X_test)

```

**2.2 K-NEAREST NEIGBORS****2.2.1 Tổng quan về K-Nearest Neighbors**

KNN (K-Nearest Neighbors) là một trong những thuật toán học có giám sát đơn giản nhất được sử dụng nhiều trong khai phá dữ liệu và học máy. Ý tưởng của thuật toán này là nó không học một điều gì từ

tập dữ liệu học (nên KNN được xếp vào loại lazy learning), mọi tính toán được thực hiện khi nó cần dự đoán nhãn của dữ liệu mới.

K-nearest neighbor có thể áp dụng được vào cả hai loại của bài toán là Classification và Regression.

Thuật toán KNN cho rằng những dữ liệu tương tự nhau sẽ tồn tại gần nhau trong một không gian, từ đó công việc của chúng ta là sẽ tìm k điểm gần với dữ liệu cần kiểm tra nhất. Việc tìm khoảng cách giữa 2 điểm cũng có nhiều công thức có thể sử dụng, tùy trường hợp mà chúng ta lựa chọn cho phù hợp. Đây là 3 cách cơ bản để tính khoảng cách 2 điểm dữ liệu x, y có k thuộc tính:

**Distance functions**

Euclidean	$\sqrt{\sum_{i=1}^k (x_i - y_i)^2}$
Manhattan	$\sum_{i=1}^k  x_i - y_i $
Minkowski	$\left( \sum_{i=1}^k ( x_i - y_i )^q \right)^{1/q}$

Hình 2.2: Các cách cơ bản để tính khoảng cách giữa 2 điểm dữ liệu

#### Ưu điểm

- Độ phức tạp tính toán của quá trình training là bằng 0.
- Việc dự đoán kết quả của dữ liệu mới rất đơn giản.
- Không cần giả sử gì về phân phối của các class.

#### Nhược điểm

- KNN rất nhạy cảm với nhiễu khi K nhỏ.

- Với K nhỏ dễ gặp nhiễu dẫn tới kết quả đưa ra không chính xác.
- Cần nhiều thời gian để thực hiện do phải tính toán khoảng cách với tất cả các đối tượng trong tập dữ liệu.

### 2.2.2 Model Implement

```

1 # import KNeighborsClassifier
2 from sklearn.neighbors import KNeighborsClassifier
3
4 # model implement
5 lr = LogisticRegression(random_state=42, solver='liblinear', multi_class='ovr')
6 lr.fit(X_train,y_train)
7
8 # predict from model
9 y_pred = lr.predict(X_test)

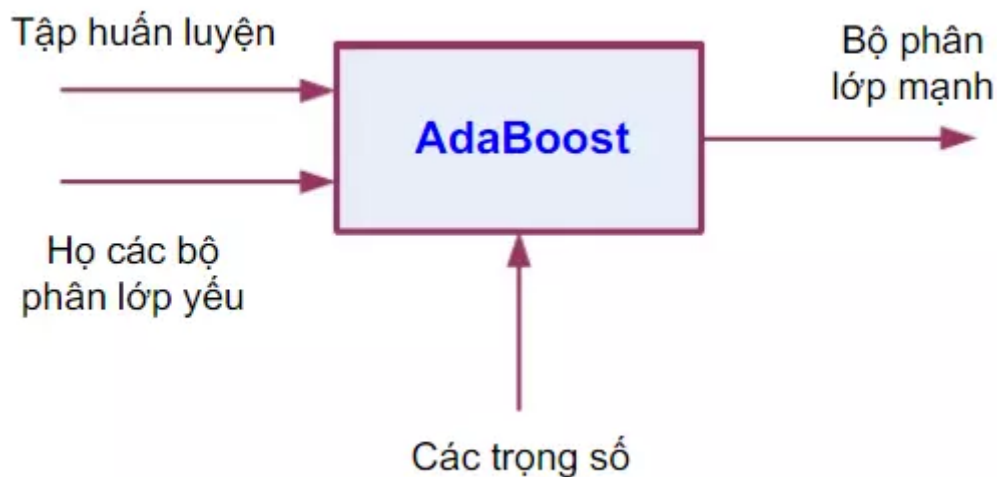
```

## 2.3 ADAPTIVE BOOSTING

### 2.3.1 Tổng quan về Adaptive Boosting

Adaptive Boost hay còn gọi là AdaBoost là một mô hình nằm trong nhánh Ensemble Learning của Machine Learning và là mô hình đầu tiên của nhánh Boosting trong Ensemble Learning. AdaBoost có tên đầy đủ là Adaptive Boosting, do Yoav Freund và Robert Schapire công bố. Với ý tưởng là sử dụng các cây quyết định để đánh trọng số cho các điểm dữ liệu, từ đó tối thiểu hóa trọng số các điểm bị phân loại sai (trọng số lớn), để tăng hiệu suất của mô hình. AdaBoost là một thuật toán học mạnh, giúp đẩy nhanh việc tạo ra một bộ phân loại mạnh bằng cách chọn các đặc trưng tốt trong một họ các bộ phân loại yếu (thường là Decision Tree) và kết hợp chúng lại tuyến tính bằng cách sử dụng các trọng số. Điều này thật sự cải thiện dần độ chính xác nhờ áp dụng hiệu quả một chuỗi các bộ phân loại yếu.

AdaBoost tiến hành train các mô hình mới dựa trên việc đánh lại trọng số cho các điểm dữ liệu hiện tại, nhằm giúp các mô hình mới có thể tập trung hơn vào các mẫu dữ liệu đang bị học sai, từ đó làm giảm giá trị loss của mô hình.



Hình 2.3: Các bước hoạt động của Adaptive Boosting

#### Ưu điểm

- AdaBoost khó dẫn đến tình trạng overfit do các thông số đầu vào không được tối ưu hóa chung
- Độ chính xác của các bộ phân loại yếu có thể được cải thiện bằng cách sử dụng AdaBoost
- Ngày nay, ngoài là một mô hình rất tốt để phân loại nhị phân, AdaBoost còn được sử dụng để phân loại text và hình ảnh

#### Nhược điểm

- AdaBoost cần một bộ data với chất lượng cao. Một bộ data với nhiều điểm nhiễu cần được xử lý trước khi sử dụng mô hình AdaBoost.
- Chỉ có thể dùng để giải quyết các vấn đề về phân loại

### 2.3.2 Model Implement

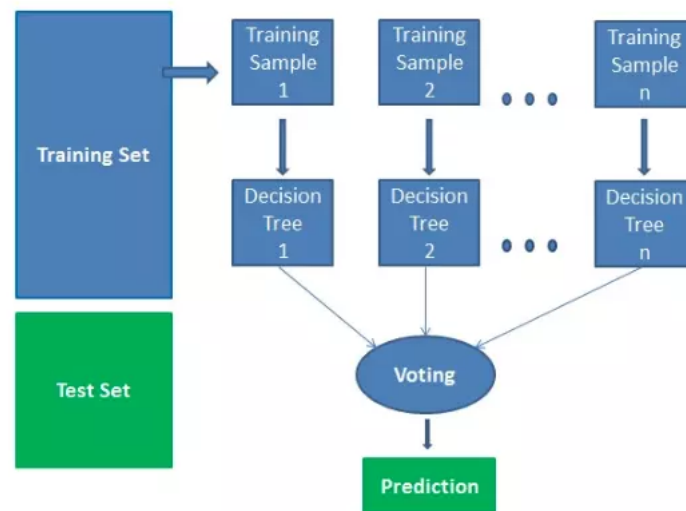
```

1 # import Adaptive Boosting
2 from sklearn.ensemble import AdaBoostClassifier
3
4 # model implement
5 ada = AdaBoostClassifier(n_estimators=100, random_state=0)
6 ada.fit(X_train, y_train)
7
8 # predict from model
9 y_pred = ada.predict(X_test)
  
```

## 2.4 RANDOM FOREST

### 2.4.1 Tổng quan về Random Forest

Random Forest là một mô hình thuộc Ensemble Learning. Random Forest được xây dựng từ nhiều cây quyết định (Decision Tree), tuy nhiên mỗi cây quyết định sẽ khác nhau (random). Sau đó kết quả được tổng hợp từ các mẫu ngẫu nhiên, được dự đoán từ mỗi cây và chọn giải pháp tốt nhất bằng cách bỏ phiếu.



Hình 2.4: Các bước hoạt động của Random Forest

Thuật toán hoạt động theo bốn bước:

- Chọn các mẫu ngẫu nhiên từ tập dữ liệu đã cho.
- Thiết lập cây quyết định cho từng mẫu và nhận kết quả dự đoán từ mỗi cây quyết định
- Bỏ phiếu cho mỗi kết quả dự đoán
- Chọn kết quả có tỉ lệ bỏ phiếu cao nhất

#### Ưu điểm

- Random Forest sử dụng nhiều decision tree và kết hợp các output của các decision tree. Điều này làm giảm overfit cho model và đồng thời giảm phương sai, từ đó tăng độ chính xác.

- Có thể được dùng cho cả các vấn đề về hồi quy và phân loại
- Có thể hoạt động tốt cả với các giá trị liên tục và phân loại
- Dễ dàng xử lý các điểm dữ liệu bị mất.
- Không yêu cầu scaling input

#### Nhược điểm

- Độ phức tạp cao do phải khởi tạo rất nhiều Decision Tree. Điều này dẫn tới việc hao phí tài nguyên
- Đi kèm với việc hao phí tài nguyên là cần nhiều thời gian để huấn luyện mô hình

#### 2.4.2 Model Implement

```

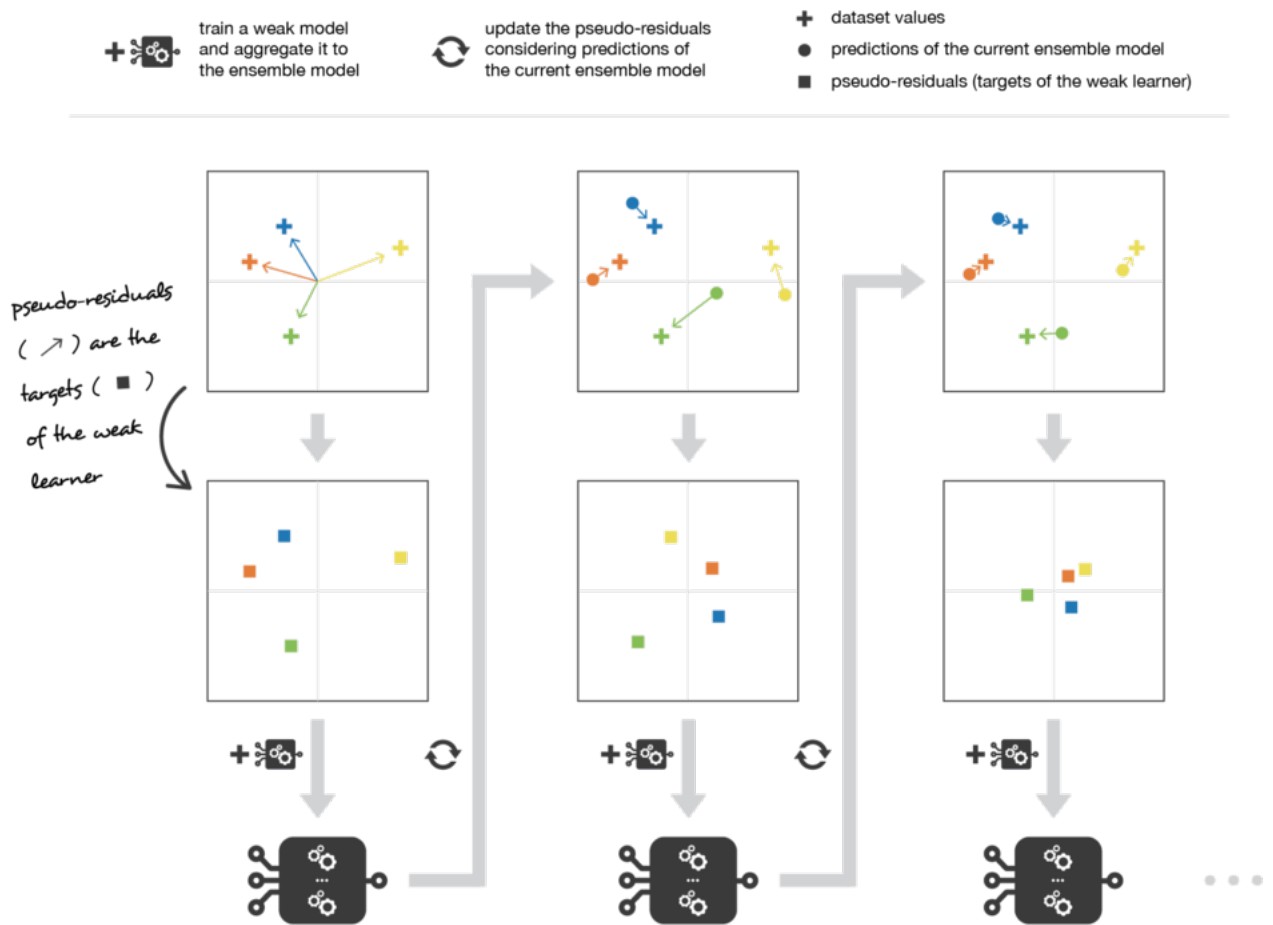
1 # import Random Forest
2 from sklearn.ensemble import RandomForestClassifier
3
4 # model implement
5 rfc = RandomForestClassifier(random_state=42)
6 rfc.fit(X_train,y_train)
7
8 # predict from model
9 y_pred = rfc.predict(X_test)

```

## 2.5 GRADIENT BOOSTING

### 2.5.1 Tổng quan về Gradient Boosting

Giống với AdaBoost, cùng thuộc nhánh boosting của Ensemble Learning, nhưng Gradient Boosting được xem là một trong những thuật toán mạnh mẽ nhất được sử dụng trong các mô hình hồi quy và phân loại. Được tạo nên từ Gradient Descent và Boosting. ‘Gradient’ nghĩa là bạn có thể sử dụng đạo hàm bậc hai hoặc đạo hàm nhiều bậc cho cùng một function. Gradient Boosting có ba phần chính, lần lượt là: additive model, loss function và một bộ phân loại yếu.



Hình 2.5: Cách hoạt động của Gradient Boosting

#### Ưu điểm

- Vì Gradient Boosting được xem như là một trong những mô hình mạnh của nhánh boosting nên độ chính xác rất cao.
- Tính linh hoạt cao, có thể tối ưu hóa trên nhiều loss function khác nhau và có nhiều các tham số để tùy chỉnh model.
- Không yêu cầu tiền xử lý dữ liệu.
- Dễ dàng xử lý các điểm dữ liệu bị mất.

#### Nhược điểm

- Gradient Boosting sẽ có thể gây ra overfit.
- Gradient Boosting cần nhiều tài nguyên để tính toán, có thể gây hao phí thời gian và bộ nhớ.



### 2.5.2 Model Implement

```
1 # import Gradient Boosting
2 from sklearn.ensemble import GradientBoostingClassifier
3
4 # model implement
5 gbc = GradientBoostingClassifier(n_estimators=100, learning_rate=0.5, max_depth=5,
6     random_state=42)
7
8 # predict from model
9 y_pred = gbc.predict(X_test)
```

## Chương 3

# DỮ LIỆU THỰC NGHIỆM

### 3.1 PHÂN TÍCH BỘ DỮ LIỆU

#### 3.1.1 Nguồn gốc data

Dữ liệu được lấy từ trang web

<https://www.kaggle.com/sakshigoyal7/credit-card-customers>

Bộ dữ liệu đã được thêm vào 2 thuộc tính mới là Home-ownership và Employees xác định khách hàng có việc làm hay không và hình thức sở hữu nhà của khách hàng đó

#### 3.1.2 Chi tiết các thuộc tính

Bộ dữ liệu bao gồm 10.000 khách hàng gồm 22 thuộc tính:

- CLIENTNUM: số id khách hàng
- Attrition\_Flag: phân loại khách hàng còn sử dụng hay đã ngưng
- Customer\_Age: Số tuổi khách hàng tính theo năm
- Gender: Giới tính khách hàng (Nam – Nữ)
- Dependent\_count: Số người phụ thuộc
- Education\_Level: Trình độ học vấn
- Marital\_Status: Tình trạng hôn nhân
- Income\_Category: Thu nhập hàng năm
- Card\_Category: Loại thẻ ngân hàng

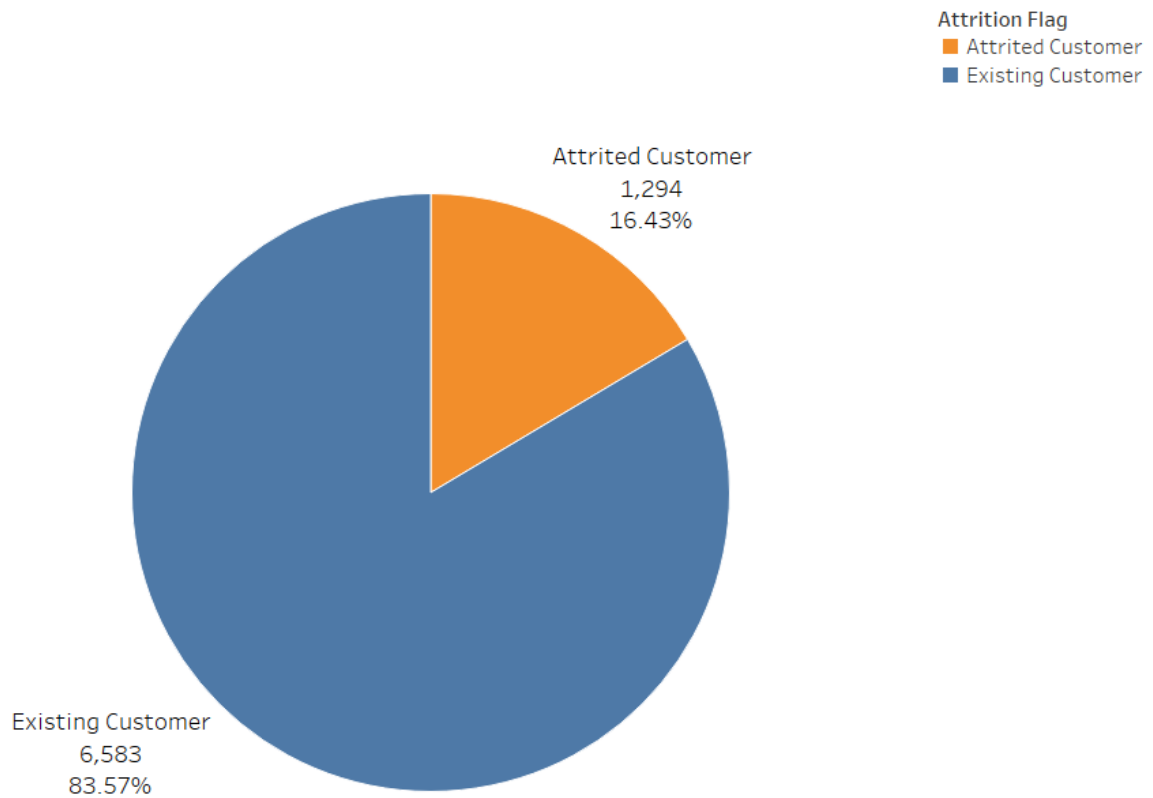
- Months\_on\_book: Thời gian dùng dịch vụ ngân hàng tính theo tháng
- Total\_Relationship\_Count: Các dịch vụ khác của ngân hàng mà khách hàng sử dụng
- Months\_Inactive\_12\_mon: Số tháng không sử dụng dịch vụ trong vòng 12 tháng
- Contacts\_Count\_12\_mon : Số lần liên hệ với ngân hàng trong vòng 12 tháng
- Credit\_Limit: Hạn mức tín dụng
- Total\_Revolving\_Bal: Số tiền mượn không trả ở mỗi cuối kì thanh toán
- Avg\_Open\_To\_Buy: Trung bình mức tối đa mà khách hàng có thể trả sử dụng thẻ ngân hàng (mức tối đa bằng hạn mức tín dụng trừ cho số dư)
- Total\_Amt\_Chng\_Q4\_Q1: thay đổi trong số tiền trong mỗi giao dịch Q4 -> Q1
- Total\_Trans\_Amt: Tổng số tiền giao dịch trong 12 tháng
- Total\_Trans\_Ct: Tổng số lần giao dịch trong 12 tháng
- Total\_Ct\_Chng\_Q4\_Q1: Thay đổi trong tổng số lần giao dịch Q4 -> Q1
- Avg\_Utilization\_Ratio: Tỷ lệ dùng dịch vụ
- Home\_ownership: Quyền sở hữu nhà của khách hàng
- Employees: Khách hàng có việc làm hay không

### 3.1.3 Trình bày bộ dữ liệu bằng Tableau

#### Churn vs Not Churn

Phân bố khách hàng theo hai nhãn đang sử dụng và ngưng sử dụng dịch vụ ngân hàng

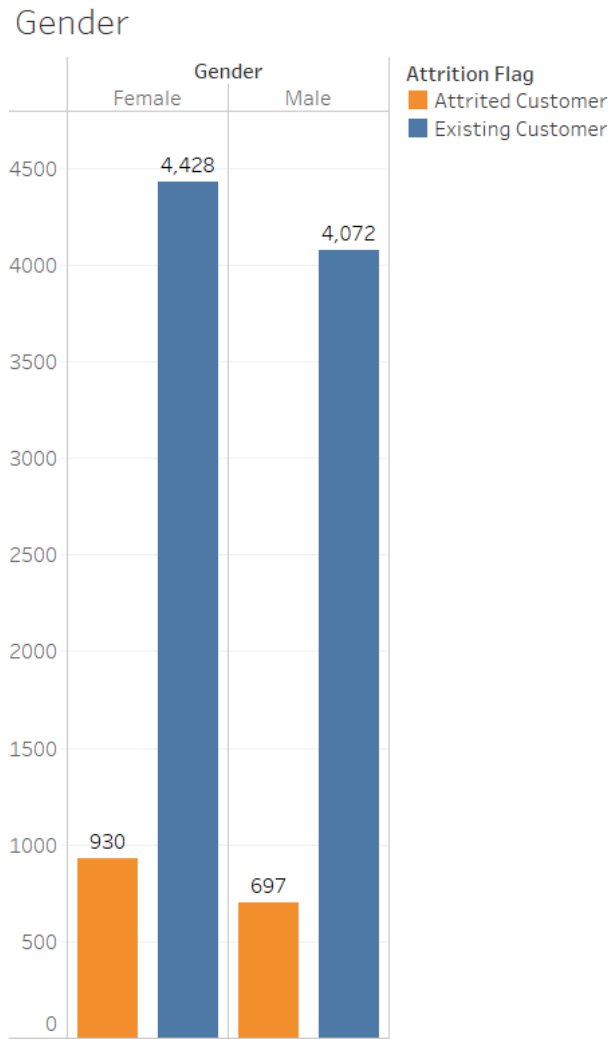
## Churn vs Not Churn



Hình 3.1: Churn vs Not Churn

**Gender**

Phân bổ khách hàng theo giới tính có thể thấy rằng số lượng khách hàng nữ bỏ đi nhiều hơn nam

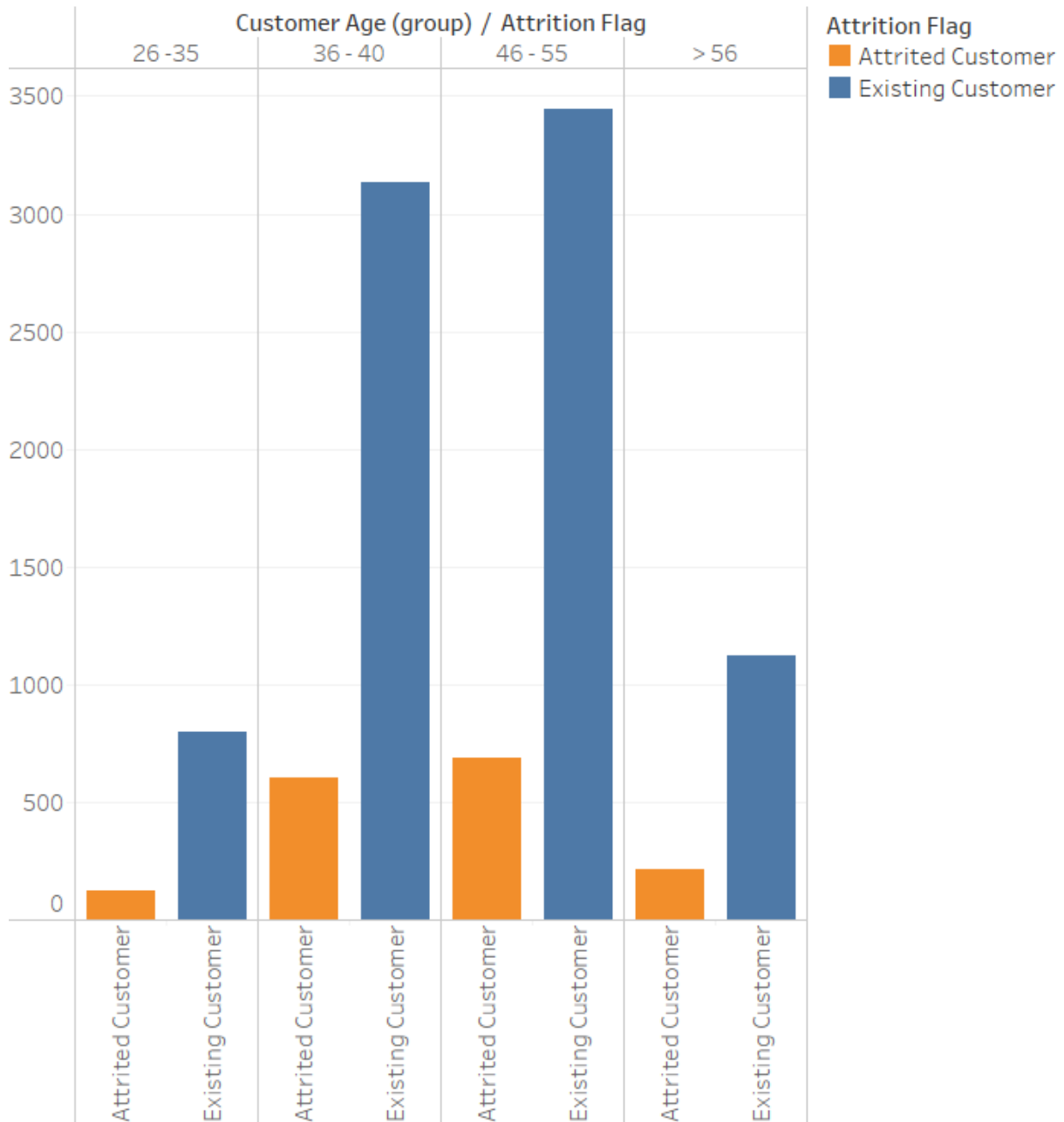


Hình 3.2: Gender

**Age group**

Phân bố khách hàng theo độ tuổi thì có thể thấy trong nhóm tuổi 36-40, 46-55 ngưng sử dụng ngân hàng nhiều nhất.

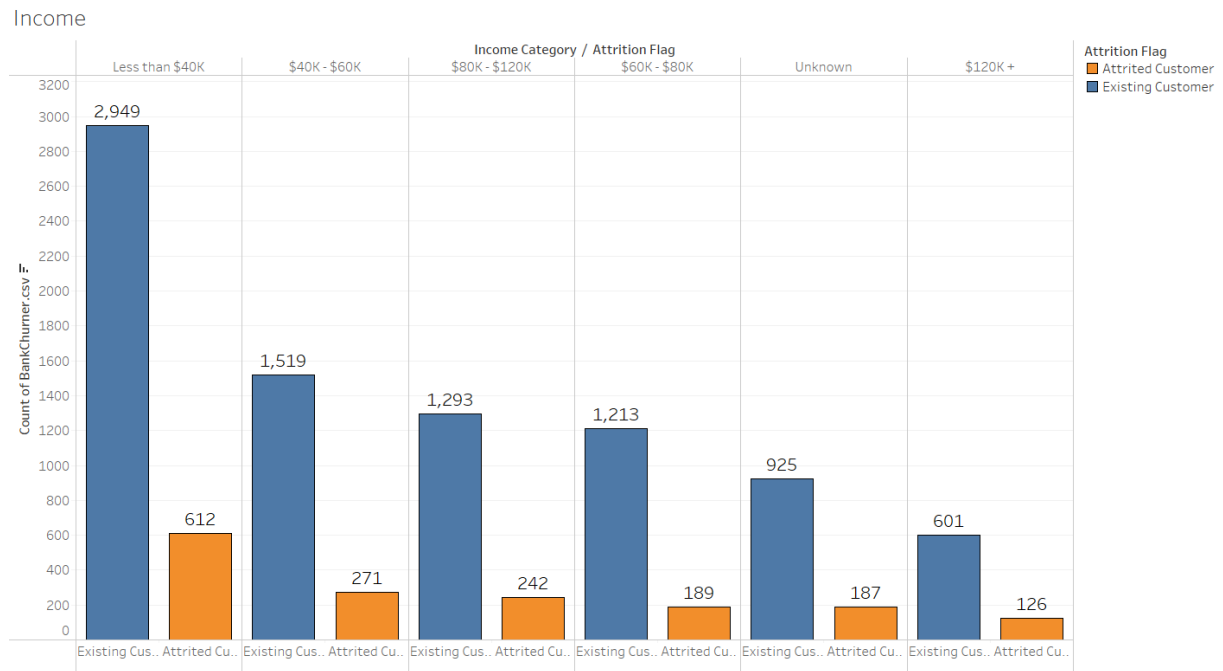
## Age Group



Hình 3.3: Age group

## Income

Phân bổ khách hàng theo thu nhập thì có thể thấy những khách hàng có thu nhập dưới 40K\$ mỗi năm ngừng sử dụng ngân hàng là nhiều nhất.

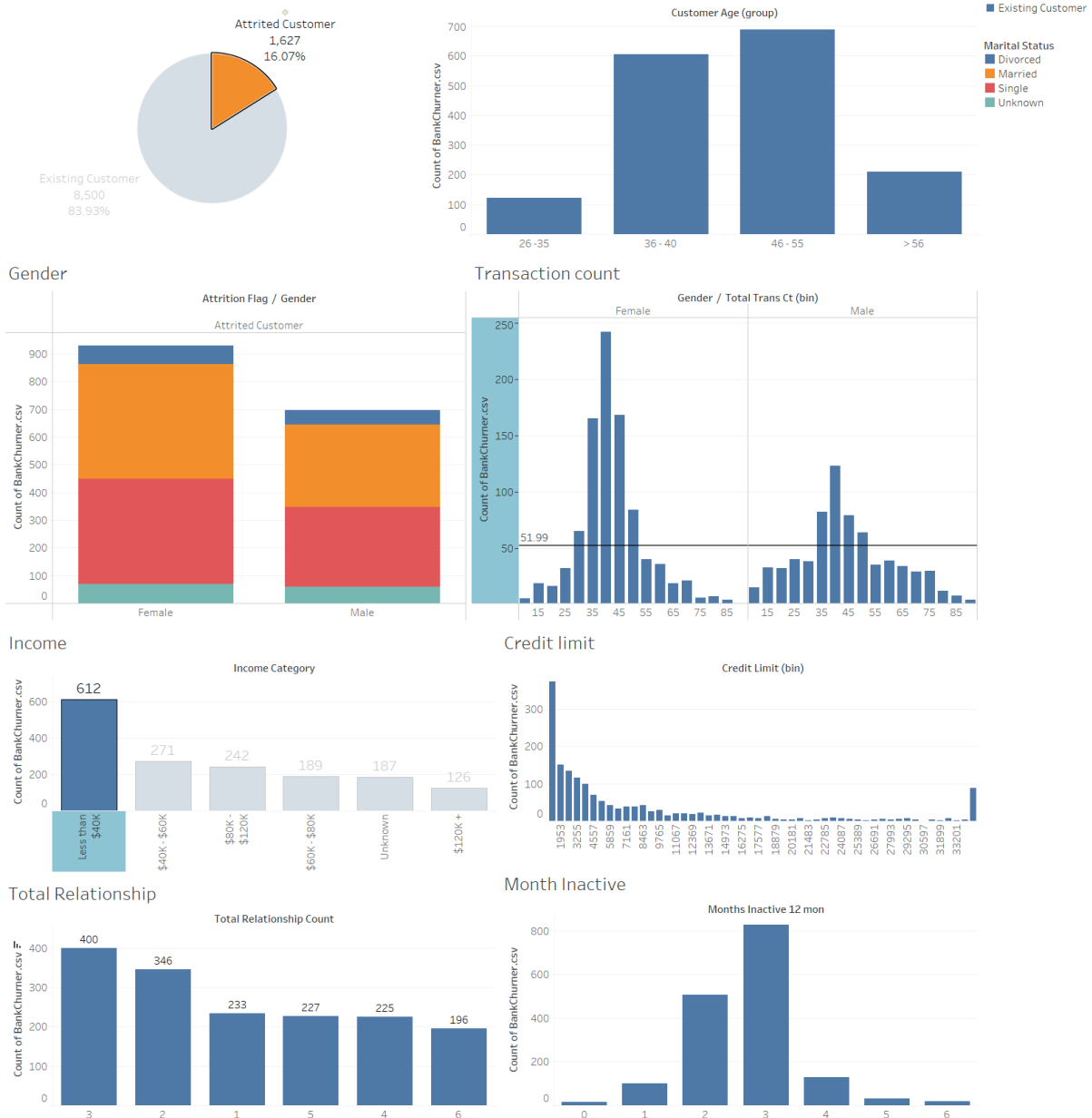


Hình 3.4: Income

### Churn Data

Khi xem xét dữ liệu của các khách hàng bỏ đi thì các khách hàng nữ trong độ tuổi 46-55 và có thu nhập dưới 40K\$ hàng năm chiếm ưu thế và đồng thời tổng số giao dịch giao động từ 35 tới 45 lần trong 1 năm.

Sheet 1



Hình 3.5: Churn Data

## Overview of Data

Tổng quan về toàn bộ dữ liệu.





Hình 3.6: Data Overview

## 3.2 TIỀN XỬ LÝ DỮ LIỆU

Ở phần tiền xử lý dữ liệu chỉ đơn giản là dùng label encoder để chuẩn hóa dữ liệu

```

1 # import library
2 from sklearn import preprocessing
3
4 le = preprocessing.LabelEncoder()
5
6 # Data preprocess
7 le.fit(df['Gender'])
8 df['Gender'] = le.transform(df['Gender'])
9 le.fit(df['Education_Level'])
10 df['Education_Level'] = le.transform(df['Education_Level'])
11 le.fit(df['Income_Category'])
12 df['Income_Category'] = le.transform(df['Income_Category'])
13 le.fit(df['Marital_Status'])
14 df['Marital_Status'] = le.transform(df['Marital_Status'])
15 le.fit(df['Card_Category'])
16 df['Card_Category'] = le.transform(df['Card_Category'])
17 le.fit(df['home_ownership'])
18 df['home_ownership'] = le.transform(df['home_ownership'])

```

```

19 le.fit(df['employees'])
20 df['employees'] = le.transform(df['employees'])
21 le.fit(df['Attrition_Flag'])
22 df['Attrition_Flag'] = le.transform(df['Attrition_Flag'])

```

CLIENTNUM	Attrition_Flag	Customer_Age	Gender	Dependent_count	Education_Level	Marital_Status	Income_Category	Card_Category	Months_on_book
768805383	Existing Customer	45	M	3	High School	Married	\$60K - \$80K	Blue	39
818770008	Existing Customer	49	F	5	Graduate	Single	Less than \$40K	Blue	44
713982108	Existing Customer	51	M	3	Graduate	Married	\$80K - \$120K	Blue	36
769911858	Existing Customer	40	F	4	High School	Unknown	Less than \$40K	Blue	34
709106358	Existing Customer	40	M	3	Uneducated	Married	\$60K - \$80K	Blue	21

Hình 3.7: 10 cột đầu tiên trước khi tiền xử lý

CLIENTNUM	Attrition_Flag	Customer_Age	Gender	Dependent_count	Education_Level	Marital_Status	Income_Category	Card_Category	Months_on_book
768805383	1	45	1	3	3	1	2	0	39
818770008	1	49	0	5	2	2	4	0	44
713982108	1	51	1	3	2	1	3	0	36
769911858	1	40	0	4	3	3	4	0	34
709106358	1	40	1	3	5	1	2	0	21

Hình 3.8: 10 cột đầu tiên sau khi tiền xử lý

```

1 # Choose feature to train
2 X_feature = {'Customer_Age', 'Gender', 'Dependent_count', 'Education_Level', '
    Marital_Status', 'Income_Category', 'Card_Category', 'Credit_Limit', '
    Months_on_book', 'Total_Relationship_Count', 'Months_Inactive_12_mon', '
    Contacts_Count_12_mon', 'Total_Revolving_Bal', 'Avg_Open_To_Buy', '
    Total_Amt_Chng_Q4_Q1', 'Total_Trans_Amt', 'Total_Trans_Ct', 'Total_Ct_Chng_Q4_Q1'
    , 'Avg_Utilization_Ratio', 'home_ownership', 'employees'}
3 X_feature = list(X_feature)
4 Y_feature = 'Attrition_Flag'
5
6 # X and y as inputs
7 X = df[X_feature]
8 y = df[Y_feature]

```

### 3.3 CÁC PHƯƠNG PHÁP ĐÁNH GIÁ GIẢI THUẬT

#### 3.3.1 Giới thiệu về các phương pháp đánh giá

Khi xây dựng các mô hình Machine Learning, chúng ta cần xác định xem mô hình đó có tốt hay không. Nhưng làm cách nào để đánh giá hiệu quả một mô hình và so sánh hiệu năng giữa các mô hình với nhau một cách khách quan nhất?

Có rất nhiều mô hình đánh giá classification. Tùy vào những bài toán khác nhau mà chúng ta sử dụng các phương pháp khác nhau. Các phương pháp thường được sử dụng là: accuracy score, confusion matrix, ROC curve, Area Under the Curve, Precision and Recall, F1 Score, etc. Tuy nhiên, trong phạm vi đề án môn học này, bọn em sẽ sử dụng các phương pháp đánh giá lần lượt là accuracy score, confusion matrix, F1 score, Precision and Recall.

Các phương pháp đánh giá này sẽ được sử dụng trên tập dữ liệu chia theo KFold với  $n\_splits=5$ . Và bộ dữ liệu gồm 2 nhãn để đánh giá. Giả sử ta gọi  $y\_pred$  là tập kết quả đầu ra được mô hình dự đoán và  $y\_test$  là tập kết quả thực.

#### 3.3.2 Accuracy Score

Accuracy score là cách đơn giản và thường được sử dụng nhất. Cách đánh giá này đơn giản là tính tỉ lệ giữa số điểm được dự đoán đúng với tổng số điểm trong tập dữ liệu kiểm thử.

```

1 import numpy as np
2
3 def acc(y_true, y_pred):
4     correct = np.sum(y_true == y_pred)
5     return float(correct)/y_true.shape[0]
6
7 y_true = np.array([0, 0, 0, 0, 1, 1, 1, 1, 0, 0])
8 y_pred = np.array([0, 1, 0, 1, 1, 1, 0, 1, 1, 0])
9 print('accuracy = ', acc(y_true, y_pred))

```

```

1 accuracy = 0.6

```

### 3.3.3 Confusion Matrix

Accuracy chỉ cho chúng ta biết về phần trăm lượng dữ liệu được phân loại đúng mà không chỉ ra được cụ thể mỗi loại được phân loại như thế nào, lớp nào được phân loại nhiều nhất, và lớp nào thường bị phân loại nhầm vào lớp khác. Ta sử dụng confusion matrix để đánh giá các giá trị trên.

Confusion matrix là một ma trận vuông với kích thước mỗi chiều bằng số lượng lớp dữ liệu. Giá trị tại hàng thứ *i*, cột thứ *j* số lượng điểm **lẽ ra thuộc vào class *i* nhưng lại được dự đoán là thuộc vào class *j***. Như vậy với phân loại nhị phân, ta sẽ có số điểm tại hàng thứ 0 cột thứ 0 là số điểm dự đoán đúng của nhãn 0, hàng thứ 0 cột thứ 1 là số điểm lẽ ra thuộc vào nhãn 0 nhưng lại được dự đoán là thuộc vào nhãn 1, tương tự với các hàng và cột còn lại.

```

1 from sklearn.metrics import confusion_matrix
2 cnf_matrix = confusion_matrix(y_true, y_pred)
3 print('Confusion matrix:')
4 print(cnf_matrix)
5 normalized_confusion_matrix = cnf_matrix/cnf_matrix.sum(axis = 1, keepdims = True)
6 print('\nConfusion matrix (with normalizatrion:)' )
7 print(normalized_confusion_matrix)

```

```

1 Confusion matrix:
2 [[3 3]
3  [1 3]]
4
5 Confusion matrix (with normalizatrion:)
6 [[0.5 0.5 ]
7  [0.25 0.75]]

```

### 3.3.4 Precision và Recall

#### True/False Positive/Negative

Ta định nghĩa True Positive(TP), False Positive(FP), True Negative(TN), False Negative(FN). Nếu xét theo confusion matrix sẽ được biểu diễn như sau:

1			Predicted		Predicted	
2			as Positive		as Negative	
3	-----		-----		-----	
4	Actual: Positive		True Positive (TP)		False Negative (FN)	
5	-----		-----		-----	
6	Actual: Negative		False Positive (FP)		True Negative (TN)	
7	-----		-----		-----	

#### Precision

Precision là tỉ lệ số điểm true positive trong số những điểm được phân loại là positive. Ta có thể khái quát bằng công thức:

$$Precision = \frac{TP}{TP + FP}$$

Precision cao đồng nghĩa với việc độ chính xác của các điểm tìm được là cao. Khi Precision = 1, mọi điểm tìm được đều thực sự là positive, tức là không có điểm negative nào lẫn vào kết quả. Tuy nhiên, Precision = 1 không đảm bảo mô hình là tốt, vì câu hỏi đặt ra là liệu mô hình đã tìm được tất cả các điểm positive hay chưa. Nếu một mô hình chỉ tìm được đúng một điểm positive mà nó chắc chắn nhất thì ta không thể gọi là một mô hình tốt.

#### Recall

Recall là tỉ lệ số điểm true positive trong số những điểm thực sự là positive. Khái quát bằng công thức như sau:

$$Recall = \frac{TP}{TP + FN}$$

Recall cao đồng nghĩa với việc True Positive cao, tức là tỉ lệ bỏ sót các điểm thực sự positive là thấp. Khi Recall = 1, mọi điểm positive đều được tìm thấy. Tuy nhiên, đại lượng này lại không đo liệu có bao nhiêu điểm negative lẫn trong đó. Nếu mô hình phân loại mọi điểm là

positive thì chắc chắn  $Recall = 1$ , tuy nhiên dễ nhận ra đây là một mô hình cực tồi.

### 3.3.5 F1-Score

F1-Score hay còn gọi là  $F_1$  score, là harmonic mean của precision và recall. Nghĩa là, F1-Score sẽ lấy cả FP và FN để tính toán. Mặc dù F1-Score khó hiểu hơn nhiều so với Accuracy nhưng F1 hữu dụng hơn nhiều so với accuracy, đặc biệt khi mà dữ liệu đầu vào bị lệch nhiều về một nhãn. Accuracy sẽ rất hiệu quả nếu FP và FN tương đương với nhau. Nếu FP và FN lệch nhau quá nhiều, F1 sẽ là sự lựa chọn nên được ưu tiên.

Công thức tính F1-Score:

$$F1 - Score = \frac{2(Precision * Recall)}{Precision + Recall}$$

## 3.4 TỔNG QUAN VỀ SMOTE

### 3.4.1 SMOTE là gì ?

SMOTE là viết tắt của Synthetic Minority Oversampling Technique. Là một kỹ thuật thống kê làm tăng số lượng dữ liệu của bộ dữ liệu một cách cân bằng. Ý tưởng của mô hình là tạo ra nhiều điểm dữ liệu từ các thuộc tính phân lớp có số lượng ít mà đầu vào cung cấp. SMOTE không thay đổi số lượng của các thuộc tính phân lớp chiếm đa số.

Các điểm dữ liệu mới không phải chỉ là bản sao chép của các điểm dữ liệu có sẵn, thay vào đó sử dụng giải thuật để tạo ra các điểm mới bằng cách kết hợp các feature của điểm gốc và feature của các điểm lân cận.

SMOTE coi toàn bộ bộ dữ liệu đầu vào làm input nhưng chỉ tăng số lượng của dữ liệu có thuộc tính phân lớp ít.

### 3.4.2 Tại sao lại cần sử dụng SMOTE ?

Trong thực tế, có nhiều bộ dữ liệu bị mất cân bằng ở các thuộc tính phân lớp. Điều này gây khó khăn cho các mô hình machine learning, đòi hỏi sự cân bằng ở các thuộc tính phân lớp để các mô hình không có sự thiên vị trong các bài toán phân loại. SMOTE là một cách để làm

tăng số lượng dữ liệu của các thuộc tính phân lớp có số lượng ít hoặc hiếm.

SMOTE trả về một bộ data mới bao gồm bộ data cũ kết hợp với lượng dữ liệu được thêm bởi SMOTE

## Chương 4

# THỰC NGHIỆM

### 4.1 IMPORT CÁC THƯ VIỆN VÀ ĐỌC DỮ LIỆU

```
1 # Import library
2 import pandas as pd
3 from sklearn.model_selection import train_test_split, cross_val_score, KFold
4 from sklearn.neighbors import KNeighborsClassifier
5 from sklearn import preprocessing
6 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
7 from sklearn.ensemble import AdaBoostClassifier, RandomForestClassifier,
8     GradientBoostingClassifier
9 import seaborn as sn
10 import matplotlib.pyplot as plt
11 import random
12 from sklearn.linear_model import LogisticRegression
13 !pip install imbalanced-learn
14 from imblearn.over_sampling import SMOTE
15
16 # Read Input from CSV
17 df = pd.read_csv('BankChurner.csv')
18 # Drop 2 unnecessary columns
19 df = df.drop(['
20     Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon
21     _Dependent_count_Education_Level_Months_Inactive_12_mon_1',
22     'Naive_Bayes_Classifier_Attrition_Flag_Card_Category_Contacts_Count_12_mon
23     _Dependent_count_Education_Level_Months_Inactive_12_mon_2'], axis = 1)
```

Tạo dữ liệu giả để đánh giá mô hình

```
1 dummy_inputs = {}
2 number_of_feature = 5
3 for feature in X:
4     dummy_inputs[feature] = []
5 for index in range(number_of_feature):
```



```

6  for feature in X:
7      dummy_inputs[feature].append(random.choice(X[feature]))
8  dummy_inputs = pd.DataFrame(dummy_inputs, index=[i for i in range(number_of_feature)
9      ])
9  dummy_inputs

```

Education_Level	Gender	Months_on_book	Credit_Limit	Total_Revolving_Bal	Card_Category	Months_Inactive_12_mon	employees	Total_Trans_Ct	Customer_Age
5	0	39	7225.0	2517	0	1	0	94	42
2	0	36	19136.0	989	0	4	1	33	51
0	1	34	29856.0	1250	0	2	1	61	52
5	0	50	2393.0	0	0	1	1	78	56
5	1	42	34516.0	0	0	3	0	55	45

Hình 4.1: 10 cột đầu tiên của dữ liệu giả dùng để đánh giá

## 4.2 ĐÁNH GIÁ TRÊN CÁC MÔ HÌNH

Các mô hình được đánh giá bằng KFold với splits = 5

```

1  kf = KFold(n_splits=5, shuffle=True, random_state=42)

```

### 4.2.1 Logistic Regression

Không sử dụng SMOTE

```

1  lr = LogisticRegression(random_state=42, solver='liblinear', multi_class='ovr')
2  ax = []
3  fig, (ax) = plt.subplots(1, 5, figsize=(30,7))
4  count = 0
5  for train_index , test_index in kf.split(X):
6      X_train , X_test = X.iloc[train_index,:],X.iloc[test_index,:]
7      y_train , y_test = y[train_index] , y[test_index]
8      lr.fit(X_train,y_train)
9      y_pred = lr.predict(X_test)
10     cfmatrix = confusion_matrix(y_pred, y_test)
11     normalize_confusion_matrix = cfmatrix / cfmatrix.sum(axis = 1, keepdims = True)
12     df_cm = pd.DataFrame(normalize_confusion_matrix, index = ["Predicted as Attrited
13     Customer", "Predicted as Existing Customer"],
14     columns = ["Attrited Customer", "Existing Customer"])
15     ax[count].set_title("Fold " + str(count + 1))
16     sn.heatmap(df_cm, cmap="YlGnBu", annot=True, ax=ax[count])
17     dummy_outputs_lr = lr.predict(dummy_inputs)
18     print("Predict dummy inputs at Fold "+str(count + 1)+": "+str(le.
19     inverse_transform(dummy_outputs_lr)))
19     print(classification_report(y_pred, y_test))

```

```

20 print("*"*60)
21 count += 1

```

## Kết quả của Logistics Regression không sử dụng SMOTE

```

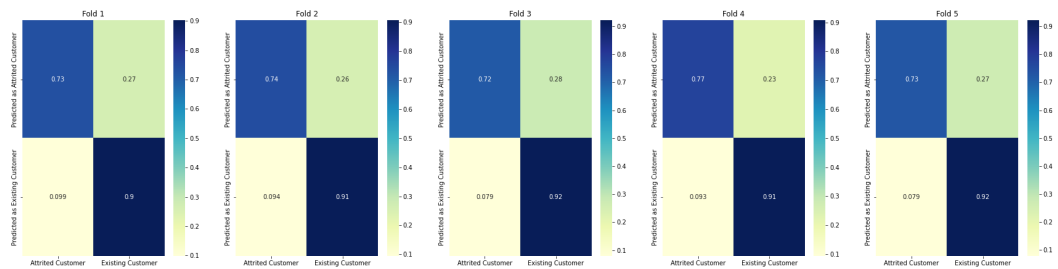
1 # Result at Fold 1
2 Predict dummy inputs at Fold 1: ['Existing Customer' 'Existing Customer' 'Existing
  Customer' 'Existing Customer' 'Existing Customer']
3           precision    recall  f1-score   support
4
5      0          0.45      0.73      0.56         200
6      1          0.97      0.90      0.93        1826
7
8      accuracy                0.88         2026
9      macro avg          0.71      0.82      0.75         2026
10     weighted avg          0.92      0.88      0.90         2026
11 # Result at Fold 2
12 Predict dummy inputs at Fold 2: ['Existing Customer' 'Existing Customer' 'Existing
  Customer' 'Existing Customer' 'Existing Customer']
13           precision    recall  f1-score   support
14
15      0          0.49      0.74      0.59         221
16      1          0.97      0.91      0.94        1805
17
18      accuracy                0.89         2026
19      macro avg          0.73      0.82      0.76         2026
20     weighted avg          0.91      0.89      0.90         2026
21 # Result at Fold 3
22 Predict dummy inputs at Fold 3: ['Existing Customer' 'Existing Customer' 'Existing
  Customer' 'Existing Customer' 'Existing Customer']
23           precision    recall  f1-score   support
24
25      0          0.54      0.72      0.62         233
26      1          0.96      0.92      0.94        1792
27
28      accuracy                0.90         2025
29      macro avg          0.75      0.82      0.78         2025
30     weighted avg          0.91      0.90      0.90         2025
31 # Result at Fold 4
32 Predict dummy inputs at Fold 4: ['Existing Customer' 'Existing Customer' 'Existing
  Customer' 'Existing Customer' 'Existing Customer']
33           precision    recall  f1-score   support
34
35      0          0.49      0.77      0.60         213
36      1          0.97      0.91      0.94        1812
37
38      accuracy                0.89         2025
39      macro avg          0.73      0.84      0.77         2025
40     weighted avg          0.92      0.89      0.90         2025
41 # Result at Fold 5

```

```

42 Predict dummy inputs at Fold 5: ['Existing Customer' 'Existing Customer' 'Existing
    Customer' 'Existing Customer' 'Existing Customer']
43           precision    recall  f1-score   support
44
45      0       0.57       0.73       0.64       251
46      1       0.96       0.92       0.94      1774
47
48   accuracy                0.90       2025
49   macro avg       0.76       0.83       0.79       2025
50   weighted avg    0.91       0.90       0.90       2025

```



Hình 4.2: Đánh giá Logistics Regression qua confusion matrix không sử dụng SMOTE

## Có sử dụng SMOTE

```

1 oversample = SMOTE()
2 X, y = oversample.fit_resample(X,y)
3
4 lr_smote = LogisticRegression(random_state=42, solver='liblinear', multi_class='ovr')
5 ax = []
6 fig, (ax) = plt.subplots(1, 5, figsize=(30,7))
7 count = 0
8 for train_index , test_index in kf.split(X):
9     X_train , X_test = X.iloc[train_index,:],X.iloc[test_index,:]
10    y_train , y_test = y[train_index] , y[test_index]
11
12    lr_smote.fit(X_train,y_train)
13    y_pred = lr_smote.predict(X_test)
14    cfmatrix = confusion_matrix(y_pred, y_test)
15    normalize_confusion_matrix = cfmatrix / cfmatrix.sum(axis = 1, keepdims = True)
16    df_cm = pd.DataFrame(normalize_confusion_matrix, index = ["Predicted as Attrited
    Customer", "Predicted as Existing Customer"],
17                        columns = ["Attrited Customer", "Existing Customer"])
18
19    ax[count].set_title("Fold " + str(count + 1))
20    sn.heatmap(df_cm, cmap="YlGnBu", annot=True, ax=ax[count])
21    dummy_outputs_lr_smote = lr_smote.predict(dummy_inputs)

```

```

22     print("Predict dummy inputs at Fold "+str(count + 1)+": "+str(1e.
        inverse_transform(dummy_outputs_lr_smote)))
23     print(classification_report(y_pred, y_test))
24     print("*"*60)
25     count += 1

```

## Kết quả của Logistics Regression có sử dụng SMOTE

```

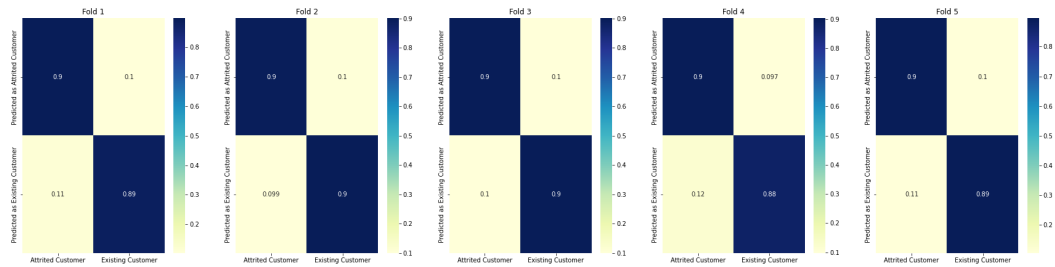
1  # Result at Fold 1
2  Predict dummy inputs at Fold 1: ['Existing Customer' 'Existing Customer' 'Existing
    Customer' 'Attrited Customer' 'Existing Customer']
3      precision    recall  f1-score   support
4
5      0          0.45      0.73      0.56         200
6      1          0.97      0.90      0.93        1826
7
8      accuracy                0.88        2026
9      macro avg          0.71      0.82      0.75        2026
10     weighted avg          0.92      0.88      0.90        2026
11 # Result at Fold 2
12 Predict dummy inputs at Fold 2: ['Existing Customer' 'Existing Customer' 'Existing
    Customer' 'Attrited Customer' 'Existing Customer']
13     precision    recall  f1-score   support
14
15     0          0.49      0.74      0.59         222
16     1          0.97      0.91      0.94        1804
17
18     accuracy                0.89        2026
19     macro avg          0.73      0.82      0.76        2026
20     weighted avg          0.91      0.89      0.90        2026
21 # Result at Fold 3
22 Predict dummy inputs at Fold 3: ['Existing Customer' 'Existing Customer' 'Existing
    Customer' 'Attrited Customer' 'Existing Customer']
23     precision    recall  f1-score   support
24
25     0          0.54      0.72      0.62         234
26     1          0.96      0.92      0.94        1791
27
28     accuracy                0.90        2025
29     macro avg          0.75      0.82      0.78        2025
30     weighted avg          0.91      0.90      0.90        2025
31 # Result at Fold 4
32 Predict dummy inputs at Fold 4: ['Existing Customer' 'Existing Customer' 'Existing
    Customer' 'Attrited Customer' 'Existing Customer']
33     precision    recall  f1-score   support
34
35     0          0.49      0.79      0.60         205
36     1          0.97      0.91      0.94        1820
37
38     accuracy                0.89        2025

```

```

39     macro avg         0.73         0.85         0.77         2025
40     weighted avg      0.93         0.89         0.90         2025
41 # Result at Fold 5
42 Predict dummy inputs at Fold 5: ['Existing Customer' 'Existing Customer' 'Existing
    Customer' 'Attrited Customer' 'Existing Customer']
43           precision      recall    f1-score      support
44
45           0           0.56         0.74         0.64         247
46           1           0.96         0.92         0.94        1778
47
48     accuracy                               0.90         2025
49     macro avg           0.76         0.83         0.79         2025
50     weighted avg        0.91         0.90         0.90         2025

```



Hình 4.3: Đánh giá Logistics Regression qua confusion matrix có sử dụng SMOTE

## 4.2.2 K-Nearest Neighbors

### Không sử dụng SMOTE

```

1 knn = KNeighborsClassifier(n_neighbors=20)
2 ax = []
3 fig, (ax) = plt.subplots(1, 5, figsize=(30,7))
4 count = 0
5 for train_index , test_index in kf.split(X):
6     X_train , X_test = X.iloc[train_index,:],X.iloc[test_index,:]
7     y_train , y_test = y[train_index] , y[test_index]
8
9     knn.fit(X_train,y_train)
10    y_pred = knn.predict(X_test)
11    cfmatrix = confusion_matrix(y_pred, y_test)
12    normalize_confusion_matrix = cfmatrix / cfmatrix.sum(axis = 1, keepdims = True)
13    df_cm = pd.DataFrame(normalize_confusion_matrix, index = ["Predicted as Attrited
    Customer", "Predicted as Existing Customer"],
14                          columns = ["Attrited Customer", "Existing Customer"])
15
16    ax[count].set_title("Fold " + str(count + 1))
17    sn.heatmap(df_cm, cmap="YlGnBu", annot=True, ax=ax[count])

```

```

18 dummy_outputs_knn = knn.predict(dummy_inputs)
19 print("Predict dummy inputs at Fold "+str(count + 1)+": "+str(le.
    inverse_transform(dummy_outputs_knn)))
20 print(classification_report(y_pred, y_test))
21 print("*"*60)
22 count += 1

```

## Kết quả của KNN không sử dụng SMOTE

```

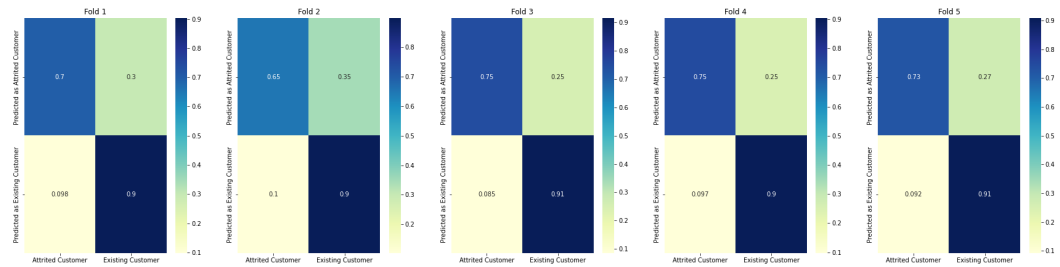
1 # Predict at Fold 1
2 Predict dummy inputs at Fold 1: ['Existing Customer' 'Existing Customer' 'Existing
    Customer' 'Attrited Customer' 'Existing Customer']
3           precision    recall  f1-score   support
4
5      0       0.46       0.70       0.55        214
6      1       0.96       0.90       0.93       1812
7
8      accuracy                0.88       2026
9      macro avg       0.71       0.80       0.74       2026
10     weighted avg       0.91       0.88       0.89       2026
11 # Predict at Fold 2
12 Predict dummy inputs at Fold 2: ['Existing Customer' 'Existing Customer' 'Existing
    Customer' 'Attrited Customer' 'Existing Customer']
13           precision    recall  f1-score   support
14
15      0       0.45       0.65       0.53        230
16      1       0.95       0.90       0.92       1796
17
18      accuracy                0.87       2026
19      macro avg       0.70       0.77       0.73       2026
20     weighted avg       0.90       0.87       0.88       2026
21 # Predict at Fold 3
22 Predict dummy inputs at Fold 3: ['Existing Customer' 'Existing Customer' 'Existing
    Customer' 'Attrited Customer' 'Existing Customer']
23           precision    recall  f1-score   support
24
25      0       0.50       0.75       0.60        207
26      1       0.97       0.91       0.94       1818
27
28      accuracy                0.90       2025
29      macro avg       0.73       0.83       0.77       2025
30     weighted avg       0.92       0.90       0.91       2025
31 # Predict at Fold 4
32 Predict dummy inputs at Fold 4: ['Existing Customer' 'Existing Customer' 'Existing
    Customer' 'Attrited Customer' 'Existing Customer']
33           precision    recall  f1-score   support
34
35      0       0.47       0.75       0.58        210
36      1       0.97       0.90       0.93       1815
37

```

```

38     accuracy                0.89    2025
39     macro avg              0.72    0.83    0.76    2025
40     weighted avg           0.92    0.89    0.90    2025
41 # Predict at Fold 5
42 Predict dummy inputs at Fold 5: ['Existing Customer' 'Existing Customer' 'Existing
    Customer' 'Attrited Customer' 'Existing Customer']
43     precision    recall    f1-score    support
44
45     0           0.48      0.73      0.58        215
46     1           0.97      0.91      0.94       1810
47
48     accuracy                0.89    2025
49     macro avg              0.72    0.82    0.76    2025
50     weighted avg           0.91    0.89    0.90    2025

```



Hình 4.4: Đánh giá KNN qua confusion matrix không sử dụng SMOTE

## Có sử dụng SMOTE

```

1 oversample = SMOTE()
2 X, y = oversample.fit_resample(X,y)
3
4 knn_smote = KNeighborsClassifier(n_neighbors=20)
5 ax = []
6 fig, (ax) = plt.subplots(1, 5, figsize=(30,7))
7 count = 0
8 for train_index , test_index in kf.split(X):
9     X_train , X_test = X.iloc[train_index,:],X.iloc[test_index,:]
10    y_train , y_test = y[train_index] , y[test_index]
11
12    knn_smote.fit(X_train,y_train)
13    y_pred = knn_smote.predict(X_test)
14    cfmatrix = confusion_matrix(y_pred, y_test)
15    normalize_confusion_matrix = cfmatrix / cfmatrix.sum(axis = 1, keepdims = True)
16    df_cm = pd.DataFrame(normalize_confusion_matrix, index = ["Predicted as Attrited
    Customer", "Predicted as Existing Customer"],
17                          columns = ["Attrited Customer", "Existing Customer"])
18

```

```

19     ax[count].set_title("Fold " + str(count + 1))
20     sn.heatmap(df_cm, cmap="YlGnBu", annot=True, ax=ax[count])
21     dummy_outputs_knn_smote = knn_smote.predict(dummy_inputs)
22     print("Predict dummy inputs at Fold "+str(count + 1)+": "+str(1e.
23         inverse_transform(dummy_outputs_knn_smote)))
24     print(classification_report(y_pred, y_test))
25     print("*"*60)
26     count += 1

```

## Kết quả của KNN có sử dụng SMOTE

```

1  # Predict at Fold 1
2  Predict dummy inputs at Fold 1: ['Attrited Customer' 'Existing Customer' 'Existing
   Customer' 'Attrited Customer' 'Existing Customer']
3      precision    recall  f1-score   support
4
5      0          0.87      0.84      0.86      1783
6      1          0.83      0.86      0.85      1617
7
8      accuracy          0.85      3400
9      macro avg          0.85      0.85      0.85      3400
10     weighted avg          0.85      0.85      0.85      3400
11 # Predict at Fold 2
12 Predict dummy inputs at Fold 2: ['Attrited Customer' 'Existing Customer' 'Existing
   Customer' 'Attrited Customer' 'Existing Customer']
13     precision    recall  f1-score   support
14
15     0          0.89      0.83      0.86      1756
16     1          0.83      0.89      0.86      1644
17
18     accuracy          0.86      3400
19     macro avg          0.86      0.86      0.86      3400
20     weighted avg          0.86      0.86      0.86      3400
21 # Predict at Fold 3
22 Predict dummy inputs at Fold 3: ['Attrited Customer' 'Existing Customer' 'Existing
   Customer' 'Attrited Customer' 'Existing Customer']
23     precision    recall  f1-score   support
24
25     0          0.87      0.83      0.85      1775
26     1          0.82      0.87      0.84      1625
27
28     accuracy          0.85      3400
29     macro avg          0.85      0.85      0.85      3400
30     weighted avg          0.85      0.85      0.85      3400
31 # Predict at Fold 4
32 Predict dummy inputs at Fold 4: ['Attrited Customer' 'Existing Customer' 'Existing
   Customer' 'Attrited Customer' 'Existing Customer']
33     precision    recall  f1-score   support
34
35     0          0.89      0.84      0.87      1830

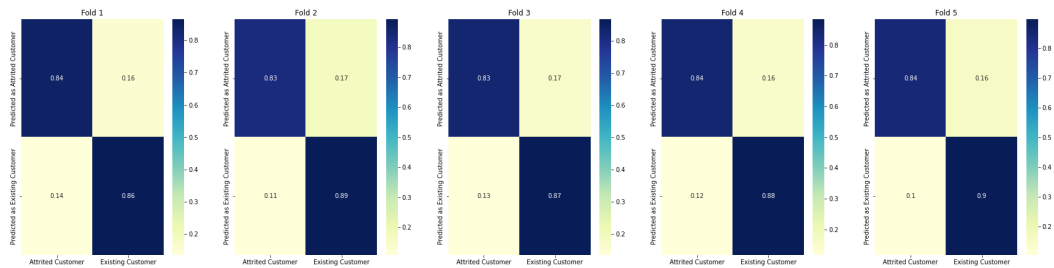
```



```

36         1         0.83         0.88         0.86         1570
37
38     accuracy                                0.86         3400
39     macro avg         0.86         0.86         0.86         3400
40     weighted avg      0.86         0.86         0.86         3400
41 # Predict at Fold 5
42 Predict dummy inputs at Fold 5: ['Attrited Customer' 'Existing Customer' 'Existing
    Customer' 'Attrited Customer' 'Existing Customer']
43         precision      recall   f1-score      support
44
45         0         0.91         0.84         0.87         1847
46         1         0.83         0.90         0.86         1553
47
48     accuracy                                0.87         3400
49     macro avg         0.87         0.87         0.87         3400
50     weighted avg      0.87         0.87         0.87         3400

```



Hình 4.5: Đánh giá KNN qua confusion matrix có sử dụng SMOTE

### 4.2.3 Adaptive Boosting

#### Không sử dụng SMOTE

```

1  ada = AdaBoostClassifier(n_estimators=100, random_state=0)
2  ax = []
3  fig, (ax) = plt.subplots(1, 5, figsize=(30,7))
4  count = 0
5  for train_index , test_index in kf.split(X):
6      X_train , X_test = X.iloc[train_index,:],X.iloc[test_index,:]
7      y_train , y_test = y[train_index] , y[test_index]
8
9      ada.fit(X_train,y_train)
10     y_pred = ada.predict(X_test)
11     cfmatrix = confusion_matrix(y_pred, y_test)
12     normalize_confusion_matrix = cfmatrix / cfmatrix.sum(axis = 1, keepdims = True)
13     df_cm = pd.DataFrame(normalize_confusion_matrix, index = ["Predicted as Attrited
        Customer", "Predicted as Existing Customer"],
14                           columns = ["Attrited Customer", "Existing Customer"])

```

```

15
16     ax[count].set_title("Fold " + str(count + 1))
17     sn.heatmap(df_cm, cmap="YlGnBu", annot=True, ax=ax[count])
18     dummy_outputs_ada = ada.predict(dummy_inputs)
19     print("Predict dummy inputs at Fold "+str(count +1)+ ":" +str(le.
20           inverse_transform(dummy_outputs_ada)))
21     print(classification_report(y_pred, y_test))
22     print("*"*60)
23     count += 1

```

## Kết quả của Adaptive Boosting không sử dụng SMOTE

```

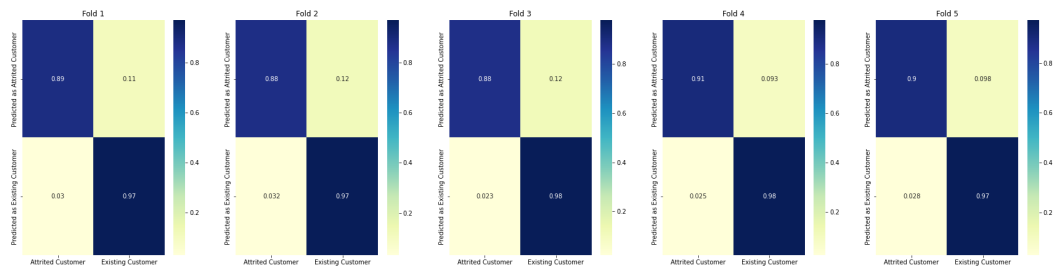
1 # Result at Fold 1
2 Predict dummy inputs at Fold 1: ['Existing Customer' 'Attrited Customer' 'Existing
3   Customer' 'Existing Customer' 'Existing Customer']
4
5           precision    recall  f1-score   support
6
7   0             0.84       0.89       0.86         310
8   1             0.98       0.97       0.97        1716
9
10  accuracy                   0.96         2026
11  macro avg           0.91       0.93       0.92         2026
12  weighted avg        0.96       0.96       0.96         2026
13 # Result at Fold 2
14 Predict dummy inputs at Fold 2: ['Existing Customer' 'Attrited Customer' 'Existing
15   Customer' 'Existing Customer' 'Existing Customer']
16
17           precision    recall  f1-score   support
18
19   0             0.84       0.88       0.86         318
20   1             0.98       0.97       0.97        1708
21
22  accuracy                   0.95         2026
23  macro avg           0.91       0.92       0.91         2026
24  weighted avg        0.95       0.95       0.95         2026
25 # Result at Fold 3
26 Predict dummy inputs at Fold 3: ['Existing Customer' 'Attrited Customer' 'Existing
27   Customer' 'Existing Customer' 'Attrited Customer']
28
29           precision    recall  f1-score   support
30
31   0             0.87       0.88       0.88         306
32   1             0.98       0.98       0.98        1719
33
34  accuracy                   0.96         2025
35  macro avg           0.92       0.93       0.93         2025
36  weighted avg        0.96       0.96       0.96         2025
37 # Result at Fold 4
38 Predict dummy inputs at Fold 4: ['Existing Customer' 'Attrited Customer' 'Existing
39   Customer' 'Existing Customer' 'Attrited Customer']
40
41           precision    recall  f1-score   support

```

```

35         0         0.87         0.91         0.89         321
36         1         0.98         0.98         0.98         1704
37
38     accuracy                                0.96         2025
39     macro avg          0.93         0.94         0.93         2025
40     weighted avg       0.97         0.96         0.96         2025
41 # Result at Fold 5
42 Predict dummy inputs at Fold 5: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Attrited Customer']
43         precision      recall    f1-score      support
44
45         0         0.85         0.90         0.88         305
46         1         0.98         0.97         0.98         1720
47
48     accuracy                                0.96         2025
49     macro avg          0.92         0.94         0.93         2025
50     weighted avg       0.96         0.96         0.96         2025

```



Hình 4.6: Đánh giá Adaptive Boosting qua confusion matrix không sử dụng SMOTE

## Có sử dụng SMOTE

```

1 oversample = SMOTE()
2 X, y = oversample.fit_resample(X,y)
3
4 ada_smote = AdaBoostClassifier(n_estimators=100, random_state=0)
5 ax = []
6 fig, (ax) = plt.subplots(1, 5, figsize=(30,7))
7 count = 0
8 for train_index , test_index in kf.split(X):
9     X_train , X_test = X.iloc[train_index,:],X.iloc[test_index,:]
10    y_train , y_test = y[train_index] , y[test_index]
11
12    ada_smote.fit(X_train,y_train)
13    y_pred = ada_smote.predict(X_test)
14    cfmatrix = confusion_matrix(y_pred, y_test)
15    normalize_confusion_matrix = cfmatrix / cfmatrix.sum(axis = 1, keepdims = True)

```

```

16 df_cm = pd.DataFrame(normalize_confusion_matrix, index = ["Predicted as Attrited Customer", "Predicted as Existing Customer"],
17                      columns = ["Attrited Customer", "Existing Customer"])
18
19 ax[count].set_title("Fold " + str(count + 1))
20 sn.heatmap(df_cm, cmap="YlGnBu", annot=True, ax=ax[count])
21 dummy_outputs_ada_smote = ada_smote.predict(dummy_inputs)
22 print("Predict dummy inputs at Fold "+str(count + 1)+ ": " + str(le.
23      inverse_transform(dummy_outputs_ada_smote)))
24 print(classification_report(y_pred, y_test))
25 print("*"*60)
26 count += 1

```

## Kết quả của Adaptive Boosting có sử dụng SMOTE

```

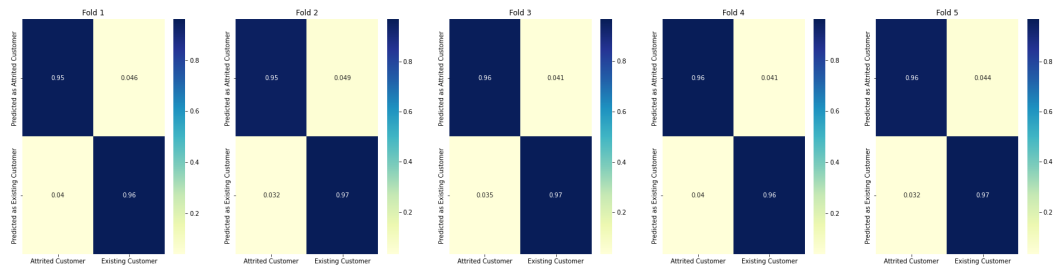
1 # Result at Fold 1
2 Predict dummy inputs at Fold 1: ['Existing Customer' 'Attrited Customer' 'Existing
   Customer' 'Existing Customer' 'Attrited Customer']
3
4      precision    recall  f1-score   support
5
6  0           0.96       0.95       0.96         1738
7  1           0.95       0.96       0.96         1662
8
9 accuracy              0.96              3400
10 macro avg           0.96           0.96           0.96           3400
11 weighted avg        0.96           0.96           0.96           3400
12 # Result at Fold 2
13 Predict dummy inputs at Fold 2: ['Existing Customer' 'Attrited Customer' 'Existing
   Customer' 'Existing Customer' 'Attrited Customer']
14
15      precision    recall  f1-score   support
16
17  0           0.97       0.95       0.96         1662
18  1           0.95       0.97       0.96         1738
19
20 accuracy              0.96              3400
21 macro avg           0.96           0.96           0.96           3400
22 weighted avg        0.96           0.96           0.96           3400
23 # Result at Fold 3
24 Predict dummy inputs at Fold 3: ['Existing Customer' 'Attrited Customer' 'Attrited
   Customer' 'Existing Customer' 'Attrited Customer']
25
26      precision    recall  f1-score   support
27
28  0           0.97       0.96       0.96         1700
29  1           0.96       0.97       0.96         1700
30
31 accuracy              0.96              3400
32 macro avg           0.96           0.96           0.96           3400
33 weighted avg        0.96           0.96           0.96           3400
34 # Result at Fold 4
35 Predict dummy inputs at Fold 4: ['Existing Customer' 'Attrited Customer' 'Attrited

```

```

Customer' 'Existing Customer' 'Attrited Customer']
33         precision    recall  f1-score   support
34
35         0          0.96      0.96      0.96     1736
36         1          0.96      0.96      0.96     1664
37
38     accuracy                0.96     3400
39     macro avg          0.96      0.96      0.96     3400
40     weighted avg       0.96      0.96      0.96     3400
41 # Result at Fold 5
42 Predict dummy inputs at Fold 5: ['Existing Customer' 'Attrited Customer' 'Existing
Customer' 'Existing Customer' 'Attrited Customer']
43         precision    recall  f1-score   support
44
45         0          0.97      0.96      0.96     1744
46         1          0.95      0.97      0.96     1656
47
48     accuracy                0.96     3400
49     macro avg          0.96      0.96      0.96     3400
50     weighted avg       0.96      0.96      0.96     3400

```



Hình 4.7: Đánh giá Adaptive Boosting qua confusion matrix có sử dụng SMOTE

## 4.2.4 Random Forest

### Không sử dụng SMOTE

```

1 rfc = RandomForestClassifier(random_state=42)
2 ax = []
3 fig, (ax) = plt.subplots(1, 5, figsize=(30,7))
4 count = 0
5 for train_index , test_index in kf.split(X):
6     X_train , X_test = X.iloc[train_index,:],X.iloc[test_index,:]
7     y_train , y_test = y[train_index] , y[test_index]
8
9     rfc.fit(X_train,y_train)
10    y_pred = rfc.predict(X_test)
11    cfmatrix = confusion_matrix(y_pred, y_test)

```

```

12     normalize_confusion_matrix = cfmatrix / cfmatrix.sum(axis = 1, keepdims = True)
13     df_cm = pd.DataFrame(normalize_confusion_matrix, index = ["Predicted as Attrited Customer", "Predicted as Existing Customer"],
14                           columns = ["Attrited Customer", "Existing Customer"])
15
16     ax[count].set_title("Fold " + str(count + 1))
17     sn.heatmap(df_cm, cmap="YlGnBu", annot=True, ax=ax[count])
18     dummy_outputs_rfc = rfc.predict(dummy_inputs)
19     print("Predict dummy inputs at Fold " + str(count + 1) + ": " + str(1e.
20           inverse_transform(dummy_outputs_rfc)))
21     print(classification_report(y_pred, y_test))
22     print("*"*60)
23     count += 1

```

## Kết quả của Random Forest không sử dụng SMOTE

```

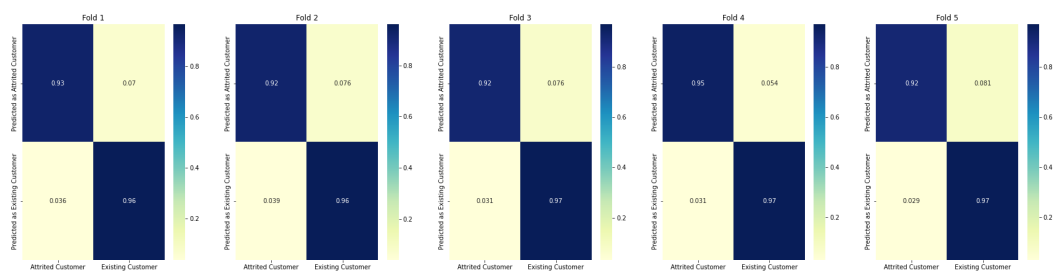
1  # Result at Fold 1
2  Predict dummy inputs at Fold 1: ['Existing Customer' 'Attrited Customer' 'Existing
   Customer' 'Existing Customer' 'Attrited Customer']
3
4
5      precision    recall  f1-score   support
6
7
8      0          0.81      0.93      0.86         284
9      1          0.99      0.96      0.98        1742
10
11
12      accuracy          0.96         2026
13      macro avg          0.90         0.95         0.92         2026
14      weighted avg          0.96         0.96         0.96         2026
15
16  # Result at Fold 2
17  Predict dummy inputs at Fold 2: ['Existing Customer' 'Attrited Customer' 'Existing
   Customer' 'Existing Customer' 'Existing Customer']
18
19
20      precision    recall  f1-score   support
21
22
23      0          0.80      0.92      0.86         288
24      1          0.99      0.96      0.97        1738
25
26
27      accuracy          0.96         2026
28      macro avg          0.89         0.94         0.91         2026
29      weighted avg          0.96         0.96         0.96         2026
30
31  # Result at Fold 3
32  Predict dummy inputs at Fold 3: ['Existing Customer' 'Attrited Customer' 'Existing
   Customer' 'Existing Customer' 'Existing Customer']
33
34
35      precision    recall  f1-score   support
36
37
38      0          0.82      0.92      0.87         276
39      1          0.99      0.97      0.98        1749
40
41
42      accuracy          0.96         2025
43      macro avg          0.91         0.95         0.92         2025
44      weighted avg          0.97         0.96         0.96         2025
45
46  # Result at Fold 4

```

```

32 Predict dummy inputs at Fold 4: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Attrited Customer']
33           precision    recall  f1-score   support
34
35         0         0.84        0.95        0.89         295
36         1         0.99        0.97        0.98        1730
37
38     accuracy                0.97        2025
39   macro avg         0.91        0.96        0.93        2025
40 weighted avg         0.97        0.97        0.97        2025
41 # Result at Fold 5
42 Predict dummy inputs at Fold 5: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Existing Customer']
43           precision    recall  f1-score   support
44
45         0         0.84        0.92        0.88         296
46         1         0.99        0.97        0.98        1729
47
48     accuracy                0.96        2025
49   macro avg         0.91        0.94        0.93        2025
50 weighted avg         0.96        0.96        0.96        2025

```



Hình 4.8: Đánh giá Random Forest qua confusion matrix không sử dụng SMOTE

## Có sử dụng SMOTE

```

1 oversample = SMOTE()
2 X, y = oversample.fit_resample(X,y)
3
4 rfc_smote = RandomForestClassifier(random_state=42)
5 ax = []
6 fig, (ax) = plt.subplots(1, 5, figsize=(30,7))
7 count = 0
8 for train_index , test_index in kf.split(X):
9     X_train , X_test = X.iloc[train_index,:],X.iloc[test_index,:]
10    y_train , y_test = y[train_index] , y[test_index]
11
12    rfc_smote.fit(X_train,y_train)

```

```

13 y_pred = rfc_smote.predict(X_test)
14 cfmatrix = confusion_matrix(y_pred, y_test)
15 normalize_confusion_matrix = cfmatrix / cfmatrix.sum(axis = 1, keepdims = True)
16 df_cm = pd.DataFrame(normalize_confusion_matrix, index = ["Predicted as Attrited Customer", "Predicted as Existing Customer"],
17                      columns = ["Attrited Customer", "Existing Customer"])
18
19 ax[count].set_title("Fold " + str(count + 1))
20 sn.heatmap(df_cm, cmap="YlGnBu", annot=True, ax=ax[count])
21 dummy_outputs_rfc_smote = rfc_smote.predict(dummy_inputs)
22 print("Predict dummy inputs at Fold "+str(count + 1)+": "+str(le.
23       inverse_transform(dummy_outputs_rfc_smote)))
24 print(classification_report(y_pred, y_test))
25 print("*"*60)
26 count += 1

```

## Kết quả của Random Forest có sử dụng SMOTE

```

1 # Result at Fold 1
2 Predict dummy inputs at Fold 1: ['Existing Customer' 'Attrited Customer' 'Existing
   Customer' 'Existing Customer' 'Attrited Customer']
3
4
5 precision    recall  f1-score   support
6
7
8      0      0.98      0.97      0.98      1747
9      1      0.97      0.98      0.98      1653
10
11 accuracy          0.98
12 macro avg          0.98
13 weighted avg       0.98
14
15 # Result at Fold 2
16 Predict dummy inputs at Fold 2: ['Existing Customer' 'Attrited Customer' 'Existing
   Customer' 'Existing Customer' 'Attrited Customer']
17
18 precision    recall  f1-score   support
19
20
21      0      0.98      0.97      0.98      1665
22      1      0.97      0.99      0.98      1735
23
24 accuracy          0.98
25 macro avg          0.98
26 weighted avg       0.98
27
28 # Result at Fold 3
29 Predict dummy inputs at Fold 3: ['Existing Customer' 'Attrited Customer' 'Existing
   Customer' 'Existing Customer' 'Attrited Customer']
30
31 precision    recall  f1-score   support
32
33
34      0      0.99      0.97      0.98      1717
35      1      0.97      0.99      0.98      1683
36
37 accuracy          0.98
38 macro avg          0.98
39 weighted avg       0.98

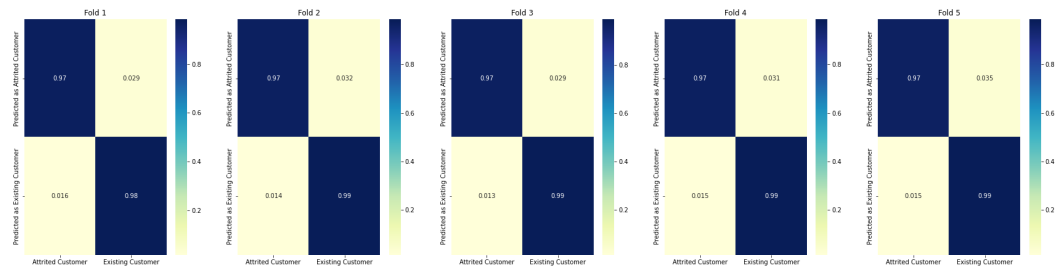
```



```

30 weighted avg      0.98      0.98      0.98      3400
31 # Result at Fold 4
32 Predict dummy inputs at Fold 4: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Attrited Customer']
33           precision      recall  f1-score      support
34
35           0           0.99      0.97      0.98      1760
36           1           0.97      0.99      0.98      1640
37
38       accuracy                        0.98      3400
39       macro avg      0.98      0.98      0.98      3400
40 weighted avg      0.98      0.98      0.98      3400
41 # Result at Fold 5
42 Predict dummy inputs at Fold 5: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Attrited Customer']
43           precision      recall  f1-score      support
44
45           0           0.99      0.97      0.98      1757
46           1           0.96      0.99      0.97      1643
47
48       accuracy                        0.97      3400
49       macro avg      0.97      0.98      0.97      3400
50 weighted avg      0.98      0.97      0.98      3400

```



Hình 4.9: Đánh giá Random Forest qua confusion matrix có sử dụng SMOTE

## 4.2.5 Gradient Boosting

### Không sử dụng SMOTE

```

1 gbc = GradientBoostingClassifier(n_estimators=100, learning_rate=0.5, max_depth=5,
    random_state=42)
2 ax = []
3 fig, (ax) = plt.subplots(1, 5, figsize=(30,7))
4 count = 0
5 for train_index , test_index in kf.split(X):
6     X_train , X_test = X.iloc[train_index,:],X.iloc[test_index,:]
7     y_train , y_test = y[train_index] , y[test_index]

```

```

8
9     gbc.fit(X_train,y_train)
10    y_pred = gbc.predict(X_test)
11    cfmatrix = confusion_matrix(y_pred, y_test)
12    normalize_confusion_matrix = cfmatrix / cfmatrix.sum(axis = 1, keepdims = True)
13    df_cm = pd.DataFrame(normalize_confusion_matrix, index = ["Predicted as Attrited
Predicted as Existing Customer"],
14                          columns = ["Attrited Customer", "Existing Customer"])
15
16    ax[count].set_title("Fold " + str(count + 1))
17    sn.heatmap(df_cm, cmap="YlGnBu", annot=True, ax=ax[count])
18
19    dummy_outputs_gbc = gbc.predict(dummy_inputs)
20    print("Predict dummy inputs at Fold "+str(count + 1)+": "+str(1e.
inverse_transform(dummy_outputs_gbc)))
21    print(classification_report(y_pred, y_test))
22    print("*"*60)
23    count += 1

```

## Kết quả của Gradient Boosting không sử dụng SMOTE

```

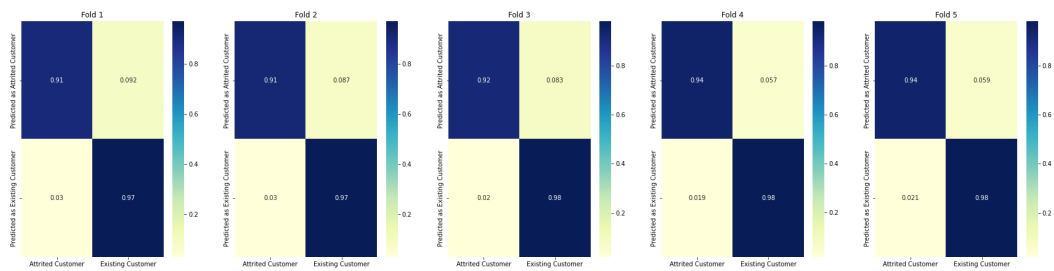
1  # Result at Fold 1
2  Predict dummy inputs at Fold 1: ['Existing Customer' 'Attrited Customer' 'Existing
Customer' 'Existing Customer' 'Attrited Customer']
3      precision    recall  f1-score   support
4
5      0          0.84      0.91      0.87         304
6      1          0.98      0.97      0.98        1722
7
8      accuracy                0.96        2026
9      macro avg          0.91      0.94      0.93        2026
10     weighted avg          0.96      0.96      0.96        2026
11 # Result at Fold 2
12 Predict dummy inputs at Fold 2: ['Existing Customer' 'Attrited Customer' 'Existing
Customer' 'Existing Customer' 'Existing Customer']
13      precision    recall  f1-score   support
14
15      0          0.85      0.91      0.88         310
16      1          0.98      0.97      0.98        1716
17
18      accuracy                0.96        2026
19      macro avg          0.92      0.94      0.93        2026
20     weighted avg          0.96      0.96      0.96        2026
21 # Result at Fold 3
22 Predict dummy inputs at Fold 3: ['Existing Customer' 'Attrited Customer' 'Existing
Customer' 'Existing Customer' 'Existing Customer']
23      precision    recall  f1-score   support
24
25      0          0.89      0.92      0.90         300
26      1          0.99      0.98      0.98        1725

```

```

27
28     accuracy                0.97    2025
29     macro avg              0.94    0.95    0.94    2025
30     weighted avg           0.97    0.97    0.97    2025
31 # Result at Fold 4
32 Predict dummy inputs at Fold 4: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Existing Customer']
33         precision    recall  f1-score   support
34
35         0          0.90      0.94      0.92         318
36         1          0.99      0.98      0.98        1707
37
38     accuracy                0.97    2025
39     macro avg              0.95    0.96    0.95    2025
40     weighted avg           0.98    0.97    0.98    2025
41 # Result at Fold 5
42 Predict dummy inputs at Fold 5: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Existing Customer']
43         precision    recall  f1-score   support
44
45         0          0.89      0.94      0.91         305
46         1          0.99      0.98      0.98        1720
47
48     accuracy                0.97    2025
49     macro avg              0.94    0.96    0.95    2025
50     weighted avg           0.97    0.97    0.97    2025

```



Hình 4.10: Đánh giá Gradient Boosting qua confusion matrix không sử dụng SMOTE

## Có sử dụng SMOTE

```

1 gbc_smote = GradientBoostingClassifier(n_estimators=100, learning_rate=0.5, max_depth
    =5, random_state=42)
2 ax = []
3 fig, (ax) = plt.subplots(1, 5, figsize=(30,7))
4 count = 0
5 for train_index , test_index in kf.split(X):
6     X_train , X_test = X.iloc[train_index,:],X.iloc[test_index,:]

```

```

7     y_train , y_test = y[train_index] , y[test_index]
8
9     gbc_smote.fit(X_train,y_train)
10    y_pred = gbc_smote.predict(X_test)
11    cfmatrix = confusion_matrix(y_pred, y_test)
12    normalize_confusion_matrix = cfmatrix / cfmatrix.sum(axis = 1, keepdims = True)
13    df_cm = pd.DataFrame(normalize_confusion_matrix, index = ["Predicted as Attrited
    Customer", "Predicted as Existing Customer"],
14                          columns = ["Attrited Customer", "Existing Customer"])
15
16    ax[count].set_title("Fold " + str(count + 1))
17    sn.heatmap(df_cm, cmap="YlGnBu", annot=True, ax=ax[count])
18    dummy_outputs_gbc_smote = gbc_smote.predict(dummy_inputs)
19    print("Predict dummy inputs at Fold "+str(count + 1)+": "+str(1e.
    inverse_transform(dummy_outputs_gbc_smote)))
20    print(classification_report(y_pred, y_test))
21    print("*"*60)
22    count += 1

```

## Kết quả của Gradient Boosting có sử dụng SMOTE

```

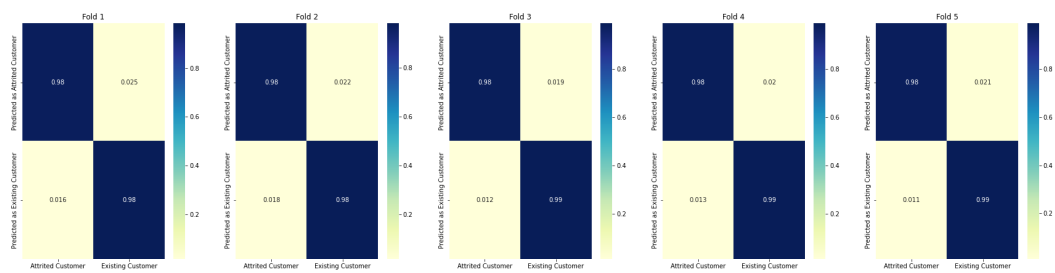
1  # Result at Fold 1
2  Predict dummy inputs at Fold 1: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Attrited Customer']
3
4      precision    recall  f1-score   support
5
6      0          0.98      0.98      0.98        1740
7      1          0.97      0.98      0.98        1660
8
9      accuracy                0.98        3400
10     macro avg          0.98      0.98      0.98        3400
11     weighted avg          0.98      0.98      0.98        3400
12 # Result at Fold 2
13 Predict dummy inputs at Fold 2: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Attrited Customer']
14
15      precision    recall  f1-score   support
16
17      0          0.98      0.98      0.98        1641
18      1          0.98      0.98      0.98        1759
19
20     accuracy                0.98        3400
21     macro avg          0.98      0.98      0.98        3400
22     weighted avg          0.98      0.98      0.98        3400
23 # Result at Fold 3
24 Predict dummy inputs at Fold 3: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Attrited Customer']
25
26      precision    recall  f1-score   support
27
28      0          0.99      0.98      0.98        1703
29      1          0.98      0.99      0.98        1697

```

```

27
28     accuracy                0.98      3400
29     macro avg              0.98      0.98      0.98      3400
30     weighted avg           0.98      0.98      0.98      3400
31 # Result at Fold 4
32 Predict dummy inputs at Fold 4: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Existing Customer']
33         precision      recall  f1-score      support
34
35         0          0.99      0.98      0.98      1744
36         1          0.98      0.99      0.98      1656
37
38     accuracy                0.98      3400
39     macro avg              0.98      0.98      0.98      3400
40     weighted avg           0.98      0.98      0.98      3400
41 # Result at Fold 5
42 Predict dummy inputs at Fold 5: ['Existing Customer' 'Attrited Customer' 'Existing
    Customer' 'Existing Customer' 'Existing Customer']
43         precision      recall  f1-score      support
44
45         0          0.99      0.98      0.98      1739
46         1          0.98      0.99      0.98      1661
47
48     accuracy                0.98      3400
49     macro avg              0.98      0.98      0.98      3400
50     weighted avg           0.98      0.98      0.98      3400

```



Hình 4.11: Đánh giá Gradient Boosting qua confusion matrix có sử dụng SMOTE

## Chương 5

# KẾT LUẬN

Trong quá trình thực hiện đề tài “Dự đoán tập khách hàng ngưng sử dụng dịch vụ ngân hàng” nhóm chúng em đã thu được các kết quả sau:

- Thêm hiểu biết về ngân hàng và các vấn đề mà các ngân hàng đang gặp phải
- Hiểu được việc rời đi của khách hàng có ảnh hưởng như thế nào đến ngân hàng
- Biết được các nhóm khách hàng nào có nguy cơ rời đi nhiều nhất.
- Hiểu rõ hơn về tính ứng dụng của học máy để dự đoán và các thuật toán học máy được áp dụng trong đề tài

Tuy nhiên vẫn còn một số khuyết điểm sau:

- Vẫn còn lúng túng trong việc trực quan hóa dữ liệu trong tableau
- Chưa hiểu sâu về kiến trúc của các thuật toán

Trong tương lai nhóm chúng em sẽ tìm cách tối ưu các mô hình đang sử dụng, cố gắng thực nghiệm dữ liệu trên các mô hình mới và phát triển thành ứng dụng triển khai được trong thực tế để có thể dự đoán được đúng thời gian khách hàng rời đi để có các chiến lược giữ lại khách hàng.

## Tài liệu tham khảo

- [1] Bui Tien, T. (2021, May 28). *Gradient Boosting - Tất tần tật về thuật toán mạnh mẽ nhất trong Machine Learning*. VIBLO ASIA. Retrieved December 9, 2021, from <https://viblo.asia/p/gradient-boosting-tat-tan-tat-ve-thuat-toan-manh-me-nhat-trong-machine-learning-YWOZrN7vZQ0>
- [2] Cao Minh, H. (2021, July 2). *AdaBoost - Bước đi đầu của Boosting*. VIBLO ASIA. Retrieved December 9, 2021, from <https://viblo.asia/p/adaboost-buoc-di-dau-cua-boosting-gAm5yrGwKdb>
- [3] CHOUDHURY, A. M. B. I. K. A. (2021, January 18). *AdaBoost Vs Gradient Boosting: A Comparison Of Leading Boosting Algorithms*. Analytics In Diamag. Retrieved December 9, 2021, from <https://analyticsindiamag.com/adaboost-vs-gradient-boosting-a-comparison-of-leading-boosting-algorithms/>
- [4] Chowdhury, M. (2021, June 26). *Churn Analytics: Data Analysis to Machine learning*. Medium. Retrieved December 6, 2021, from <https://medium.com/@mchowdhuryca/churn-analytics-from-data-analysis-to-machine-learning-95854d102ed6>
- [5] DHANYA, S. H. R. E. E. (2021, June 1). *AdaBoost : A Brief Introduction to Ensemble learning*. Analytics Vidhya. Retrieved December 10, 2021, from <https://www.analyticsvidhya.com/blog/2021/06/adaboost-a-brief-introduction-to-ensemble-learning/>
- [6] Konstantin, T. (2021, May 1). *Bank Churn Data Exploration And Churn Prediction*. Kaggle. Retrieved November 5, 2021, from <https://www.kaggle.com/thomaskonstantin/bank-churn-data-exploration-and-churn-prediction>
- [7] Kumar, N. (2019, February 23). *Advantages and Disadvantages of Random Forest Algorithm in Machine Learning*. The Professionals Point. Retrieved December 11, 2021, from

- <https://theprofessionalspoint.blogspot.com/2019/02/advantages-and-disadvantages-of-random.html>
- [8] Microsoft. (2021, August 23). *SMOTE*. Retrieved December 7, 2021, from <https://docs.microsoft.com/en-us/azure/machine-learning/studio-module-reference/smote>
  - [9] Pham Minh, H. (2021, January 17). *Ensemble learning và các biến thể (P1)*. VIBLO ASIA. Retrieved December 9, 2021, from <https://viblo.asia/p/ensemble-learning-va-cac-bien-the-p1-WAyK80AkKxX>
  - [10] Pham Minh, P. (2021, May 20). *Giải thuật Adaboost ứng dụng trong nhận dạng biến số xe(Data Mining)*. VIBLO ASIA. Retrieved December 9, 2021, from <https://viblo.asia/p/giai-thuat-adaboost-ung-dung-trong-nhan-dang-bien-so-xedata-mining-07LKXxbPKV4>
  - [11] R. (2021, November 15). *Know your customers to know your churners*. Kaggle. Retrieved November 6, 2021, from <https://www.kaggle.com/renmellofrey/know-your-customers-to-know-your-churners>
  - [12] *Random forest Algorithm in Machine learning: An Overview*. (2020, February 19). Great Learning. Retrieved December 10, 2021, from <https://www.mygreatlearning.com/blog/random-forest-algorithm/>
  - [13] Tuong, H. (2020, August 15). *Thuật Toán K-Nearest Neighbors (KNN) Siêu Cơ Bản*. CODELEARN. Retrieved December 7, 2021, from <https://codelearn.io/sharing/thuat-toan-k-nearest-neighbors-knn>
  - [14] Vu Huu, T. (2017a, January ). *K-nearest neighbors*. Machine Learning Cơ Bản. Retrieved December 7, 2021, from <https://machinelearningcoban.com/2017/01/08/knn/>
  - [15] Vu Huu, T. (2017b, January 27). *Logistic Regression*. Machine Learning Cơ Bản. Retrieved December 8, 2021, from <https://machinelearningcoban.com/2017/01/27/logisticregression/>
  - [16] Vu Huu, T. (2018, January 3). *Các phương pháp đánh giá một hệ thống phân lớp*. Machine Learning Cơ Bản. Retrieved December 11, 2021, from <https://machinelearningcoban.com/2017/08/31/evaluation/>