

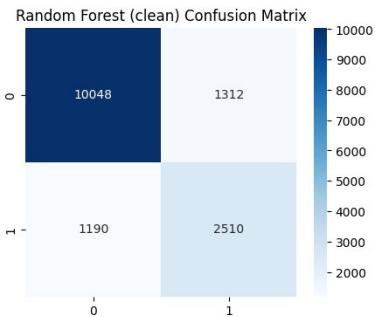
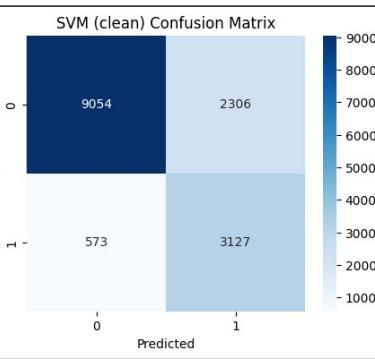
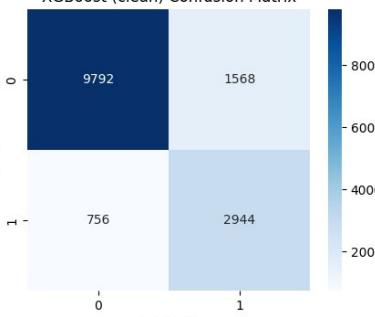
# Bài thực hành buổi 2

## Bài 1

Kết quả chạy 3 mô hình đối với dữ liệu chưa được xử lý

	Accuracy	Precision	Recall	F1	Confusion Matrix																
Gradient Boosting	0.8683	0.7853	0.6386	0.7044	<p>Gradient Boosting(raw) Confusion Matrix</p> <table border="1"><caption>Gradient Boosting(raw) Confusion Matrix</caption><thead><tr><th colspan="2">Actual</th><th colspan="2">Predicted</th></tr><tr><th>0</th><th>1</th><th>0</th><th>1</th></tr></thead><tbody><tr><th>0</th><td>10714</td><td>646</td><td></td></tr><tr><th>1</th><td>1337</td><td>2363</td><td></td></tr></tbody></table>	Actual		Predicted		0	1	0	1	0	10714	646		1	1337	2363	
Actual		Predicted																			
0	1	0	1																		
0	10714	646																			
1	1337	2363																			
SVM	0.7908	0.9678	0.1538	0.2654	<p>SVM(raw) Confusion Matrix</p> <table border="1"><caption>SVM(raw) Confusion Matrix</caption><thead><tr><th colspan="2">Actual</th><th colspan="2">Predicted</th></tr><tr><th>0</th><th>1</th><th>0</th><th>1</th></tr></thead><tbody><tr><th>0</th><td>10651</td><td>709</td><td></td></tr><tr><th>1</th><td>1482</td><td>2218</td><td></td></tr></tbody></table>	Actual		Predicted		0	1	0	1	0	10651	709		1	1482	2218	
Actual		Predicted																			
0	1	0	1																		
0	10651	709																			
1	1482	2218																			
XGBoost	0.8012	0.5924	0.6116	0.6019	<p>XGB(raw) Confusion Matrix</p> <table border="1"><caption>XGB(raw) Confusion Matrix</caption><thead><tr><th colspan="2">Actual</th><th colspan="2">Predicted</th></tr><tr><th>0</th><th>1</th><th>0</th><th>1</th></tr></thead><tbody><tr><th>0</th><td>10660</td><td>700</td><td></td></tr><tr><th>1</th><td>1246</td><td>2454</td><td></td></tr></tbody></table>	Actual		Predicted		0	1	0	1	0	10660	700		1	1246	2454	
Actual		Predicted																			
0	1	0	1																		
0	10660	700																			
1	1246	2454																			

### Kết quả chạy 3 mô hình đối với dữ liệu đã được xử lý

	Accuracy	Precision	Recall	F1	Confusion Matrix												
<b>Gradient Boosting</b>	0.8562	0.6984	0.73	0.7139	 <p>Random Forest (clean) Confusion Matrix</p> <table border="1"> <thead> <tr> <th colspan="2"></th> <th>Actual</th> <th>0</th> <th>1</th> </tr> <tr> <th rowspan="2">Actual</th> <th>0</th> <td>10048</td> <td>1312</td> </tr> </thead> <tbody> <tr> <th>1</th> <td>1190</td> <td>2510</td> </tr> </tbody> </table>			Actual	0	1	Actual	0	10048	1312	1	1190	2510
		Actual	0	1													
Actual	0	10048	1312														
	1	1190	2510														
<b>SVM</b>	0.8373	0.6443	0.7538	0.6947	 <p>SVM (clean) Confusion Matrix</p> <table border="1"> <thead> <tr> <th colspan="2"></th> <th>Actual</th> <th>0</th> <th>1</th> </tr> <tr> <th rowspan="2">Actual</th> <th>0</th> <td>9054</td> <td>2306</td> </tr> </thead> <tbody> <tr> <th>1</th> <td>573</td> <td>3127</td> </tr> </tbody> </table>			Actual	0	1	Actual	0	9054	2306	1	573	3127
		Actual	0	1													
Actual	0	9054	2306														
	1	573	3127														
<b>XGBoost</b>	0.8132	0.6199	0.6197	0.6198	 <p>XGBoost (clean) Confusion Matrix</p> <table border="1"> <thead> <tr> <th colspan="2"></th> <th>Actual</th> <th>0</th> <th>1</th> </tr> <tr> <th rowspan="2">Actual</th> <th>0</th> <td>9792</td> <td>1568</td> </tr> </thead> <tbody> <tr> <th>1</th> <td>756</td> <td>2944</td> </tr> </tbody> </table>			Actual	0	1	Actual	0	9792	1568	1	756	2944
		Actual	0	1													
Actual	0	9792	1568														
	1	756	2944														

### NHẬN XÉT:

- Tiền xử lý làm tăng Recall và cải thiện F1 ở cả ba mô hình: so với dữ liệu chưa xử lý, Gradient Boosting tăng F1 từ 0.7044 lên 0.7139, SVM tăng mạnh từ 0.2654 lên 0.6947, XGBoost tăng nhẹ từ 0.6019 lên 0.6198

Do tiền xử lý giúp làm rõ ranh giới giữa các lớp, tăng tính nhất quán của đặc trưng cũng như làm cân bằng hơn giữa precision và recall nên các mô hình nhạy cảm với scale hoặc noise (như SVM) cải thiện rất nhiều; mô hình mạnh với dữ liệu thô (XGBoost) chỉ tăng nhẹ do trước đó nó đã học nhiều đặc trưng, còn Gradient Boosting thể hiện ổn định và vẫn cải thiện nhẹ

- Gradient Boosting đạt giá trị cao nhất ở tất cả các chỉ số sau tiền xử lý (Accuracy, Precision, Recall và F1-score) => Đây là mô hình hoạt động tốt nhất trên tập dữ liệu của đề tài

Lý do là Gradient Boosting học theo cơ chế tăng dần độ chính xác: mỗi cây kế tiếp được xây dựng để sửa lỗi của cây trước, nên mô hình vừa đủ linh hoạt để học quan hệ giữa các đặc trưng, nhưng không quá phức tạp đến mức học theo nhiễu. Đồng thời mô hình cũng không phụ thuộc mạnh vào việc chuẩn hóa dữ liệu nên duy trì hiệu năng ổn định giữa hai phiên bản dữ liệu. Nhờ sự cân bằng giữa độ phức tạp và khả năng tổng quát này, Gradient Boosting thể hiện tốt hơn các mô hình còn lại trong bài toán.

## **Bài 2:**

### **Nhận xét:**

Từ kết quả vừa chạy, có thể thấy rằng accuracy trung bình của dữ liệu gốc là 73.09%, trong khi dữ liệu tăng cường đạt 75.89%, nghĩa là augmentation giúp cải thiện độ chính xác khoảng 2.8%. Điều này cho thấy việc tăng cường dữ liệu giúp mô hình học tổng quát hơn và giảm overfitting. Về thời gian huấn luyện, dữ liệu gốc mất trung bình 250.8 giây, còn dữ liệu tăng cường mất 585.78 giây, tăng đáng kể do số lượng ảnh lớn hơn và các phép biến đổi tốn thời gian tính toán.

ξ Augmentation giúp tăng độ chính xác nhưng thời gian huấn luyện sẽ lâu hơn