

In this task, each step was carefully designed to optimize model performance and address specific challenges:

1. **Stratified Splitting:**

Since the dataset was imbalanced, I used stratified splitting to preserve class distribution across the training and testing sets. This approach ensures fair evaluation across all classes and mitigates bias towards the majority class.

2. **Bayesian Optimization:**

To select the best hyperparameters efficiently, I used Bayesian optimization instead of traditional grid or random search. Bayesian optimization is faster and more effective as it utilizes prior evaluations to focus on the most promising areas of the hyperparameter space.

3. **Testing Multiple Models:**

I evaluated four different models: SVM, Logistic Regression, Random Forest, and XGBoost. Each model was tested under two scenarios: with and without PCA. This allowed for a comprehensive comparison to identify the best-performing model.

4. **PCA (Principal Component Analysis):**

PCA was applied to reduce the dataset's dimensionality, simplifying the training process and reducing the risk of overfitting. This also improved computational efficiency. Based on the results, XGBoost with 43 PCA components outperformed the other models.

5. **Random State Fine-tuning:**

After determining that XGBoost with PCA provided the best performance, I further fine-tuned the model by experimenting with different random states for both data splitting and model initialization. This helped achieve more consistent and reliable performance.

6. **Model Performance:**

The best configuration achieved an accuracy of 80% and an F1 score of 0.79. These metrics indicate a strong overall performance, particularly given the imbalanced nature of the dataset.

7. **Calibration:**

Finally, I calibrated the model to improve the reliability of its predicted probabilities. While this slightly reduced the accuracy, it made the predicted probabilities more trustworthy. Calibration ensures that the model's probability estimates align better with actual outcomes, which is crucial for tasks that rely on probabilistic decision-making.