<span style="color:green">**Flaw Introduced in the Random Forest Model:**</span>

<span style="color:green">**Flaw: Imbalanced Training Dataset**</span>

- **Description**:
    - The training dataset was deliberately manipulated to create a significant imbalance. The majority class (fully_paid = 1) was oversampled to comprise **95% of the training data**, while the minority class (fully_paid = 0) was undersampled to only **5% of the training data**. This imbalance biases the model toward predicting the majority class.

---

**Impact on the Model and Results:**

**1. Predictive Bias:**

- The model becomes **overly reliant on the majority class (fully_paid = 1)** during training.

- As a result, it struggles to correctly identify the minority class (fully_paid = 0), leading to:
    - **High recall for the majority class.**
    - **Very low precision and recall for the minority class.**

**2. Misleading Metrics:**

- Despite the imbalance, the **ROC AUC score remains artificially high (0.9571)** because it measures the model's ability to rank predictions, which is not directly impacted by imbalance in many cases.

- **Accuracy** is extremely low (0.5017), reflecting the model's inability to generalize due to the flawed training process.

- The **classification report** shows the stark contrast:
    - Precision for the minority class (0) is **1.00** (indicating no false positives) because the model predicts almost everything as 1 in the test set.
    - Recall for the minority class is **0.00**, meaning the model fails to identify any minority instances.

**3. Feature Importance Skew:**

- Features most relevant to the minority class might appear less important because the model does not encounter enough diverse data to learn their significance.

- This misrepresentation affects interpretability and downstream decision-making.

**4. Generalization Failure:**

- The model performs poorly when exposed to real-world data where the class distribution might differ or when minority instances are critical (e.g., identifying risky loans).

- In real applications, this could result in **financial risk** (approving bad loans) or **unfair decisions** (denying eligible loans).

**Real-World Risks:**

- **Loan Approvals**: The model might approve loans for high-risk individuals, leading to increased defaults.

- **Decision-Making Bias**: Critical features for identifying default risks might be undervalued.

- **Compliance Issues**: Regulatory frameworks often mandate fairness in automated decision-making, which this model violates by exhibiting bias.

---

**Key Takeaway:**

This flaw highlights how training data imbalance can drastically affect a model's performance, leading to skewed predictions, poor generalization, and potentially severe real-world consequences in financial applications.

---

In this **XGBoost model with a flaw**, the introduced issue is **incorrect labeling in the target variable**:

**Flaw Description**

- **Flipping Labels**:

  - A portion (10%) of the training labels (y_train) were randomly flipped, changing some 1 (fully paid) labels to 0 (charged off) and vice versa.

  - This deliberate corruption of the training data introduces noise into the learning process, making the target variable inconsistent.

**Impact of the Flaw**

1. **Model Confusion**:

   - The model struggles to find accurate patterns in the training data due to conflicting signals caused by incorrect labels.

   - This affects the model's ability to generalize to unseen data during testing.

2. **Degraded Accuracy**:

   - The training performance might seem adequate due to overfitting on the noisy labels, but the model's true predictive power on clean test data is compromised.

   - Even if the ROC AUC appears reasonable, the model's predictions may not be trustworthy in real-world applications.

3. **Misleading Feature Importance**:

- Features correlated with the flipped labels may gain or lose importance inappropriately, leading to incorrect interpretations of the drivers behind loan repayment or default.

4. **Operational Risks**:

   - Deploying such a model could result in poor decision-making, such as granting loans to high-risk individuals or denying loans to low-risk applicants.

   - The flaw undermines the reliability of the model in production settings.

**Indicators of the Flaw**

- **High AUC but Low Precision/Recall**:

  - The model might still achieve a good AUC (due to the test set not being flawed), but classification metrics such as precision, recall, or F1-score may indicate inconsistencies.

- **Performance Gap**:

  - A significant gap between the training and testing performance could suggest issues with the data used during training.

By introducing this flaw, we demonstrate how critical clean and accurate labeling is to building reliable machine learning models.