

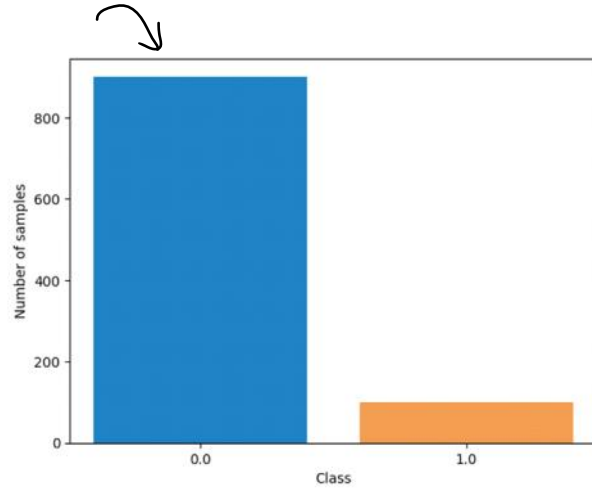
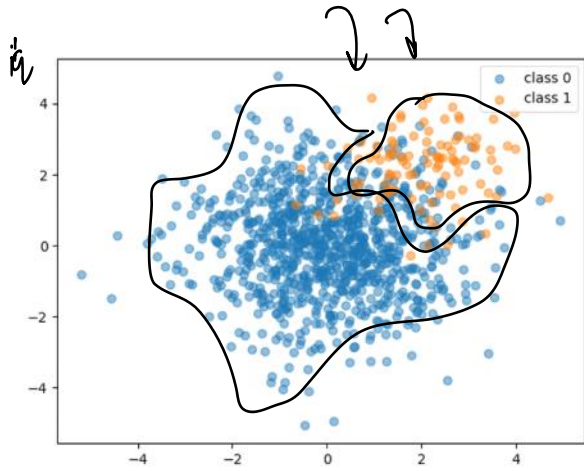
What is Imbalanced Data  $\rightarrow$  classification  $\rightarrow$  cgp (classification)

24 April 2024 07:32

Imbalanced data refers to a situation in a dataset where the classes or categories are not equally represented. It can happen in both binary classification setup as well as multi-class classification setup.

$$\begin{array}{r} 450 - 0 \\ \hline 50 - 1 \end{array}$$

↑  
imbal

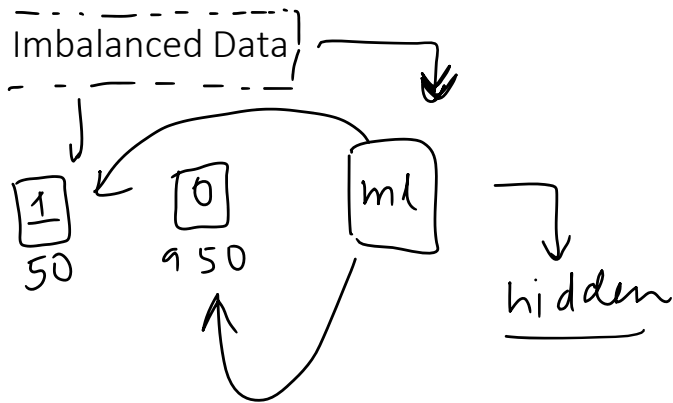


cgp

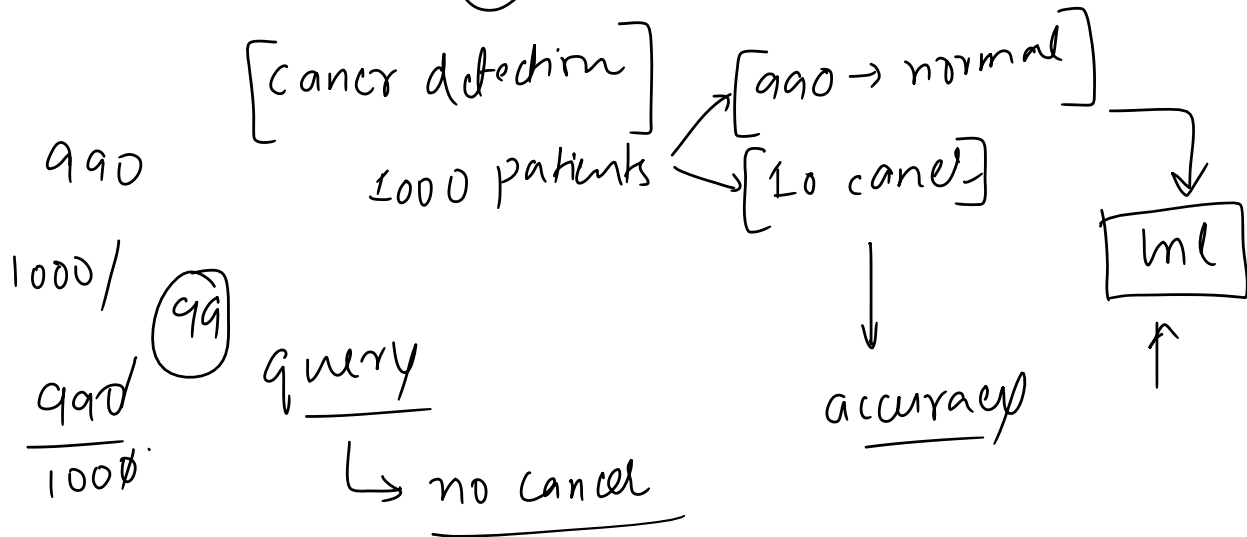
$$\begin{bmatrix} 0 & 1 & 2 \\ 800 & 100 & 100 \end{bmatrix}$$

# Problem with Imbalanced Data

24 April 2024 07:37



accuracy → misleading



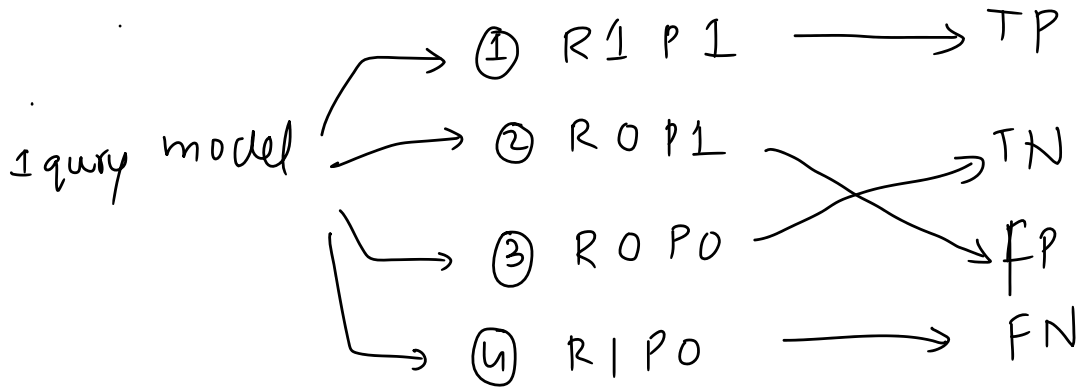
# Cancer detection

TP

TN

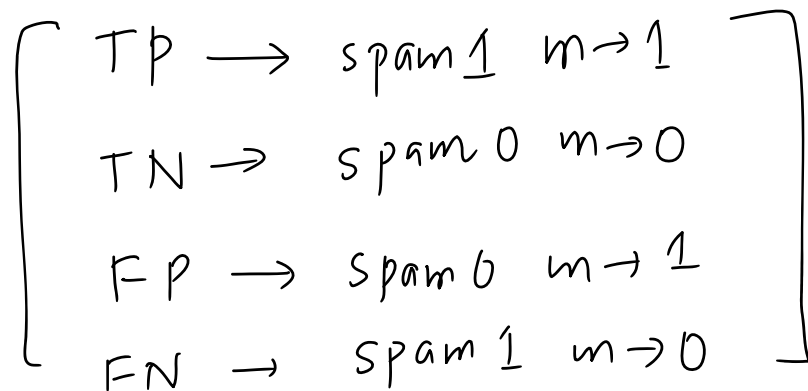
FP

FN

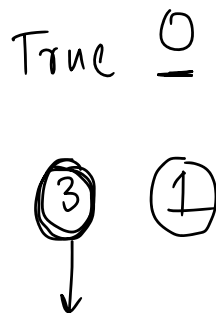
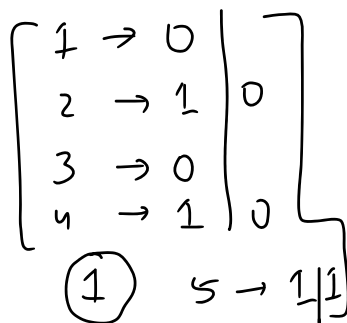


## Spam detection

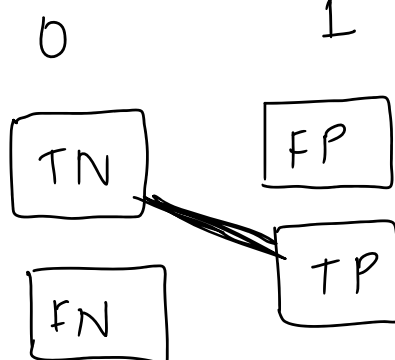
0 not spam 1 spam



5 points → model



Pred



confusion matrix

Accuracy

Precision

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

$$= \frac{1}{3}$$

Precision (low)

FP

$$\text{accuracy} = \frac{TP + TN}{TN + TP + FP + FN}$$

FP

FN

FP

spam

spam 1  
m 0

spam 0  
m 1

FN

recall

FP

problem

① →

Recall

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

Recall low

→ FN ↑

50 50 99 1  
51 49 [61 39]  
60 40

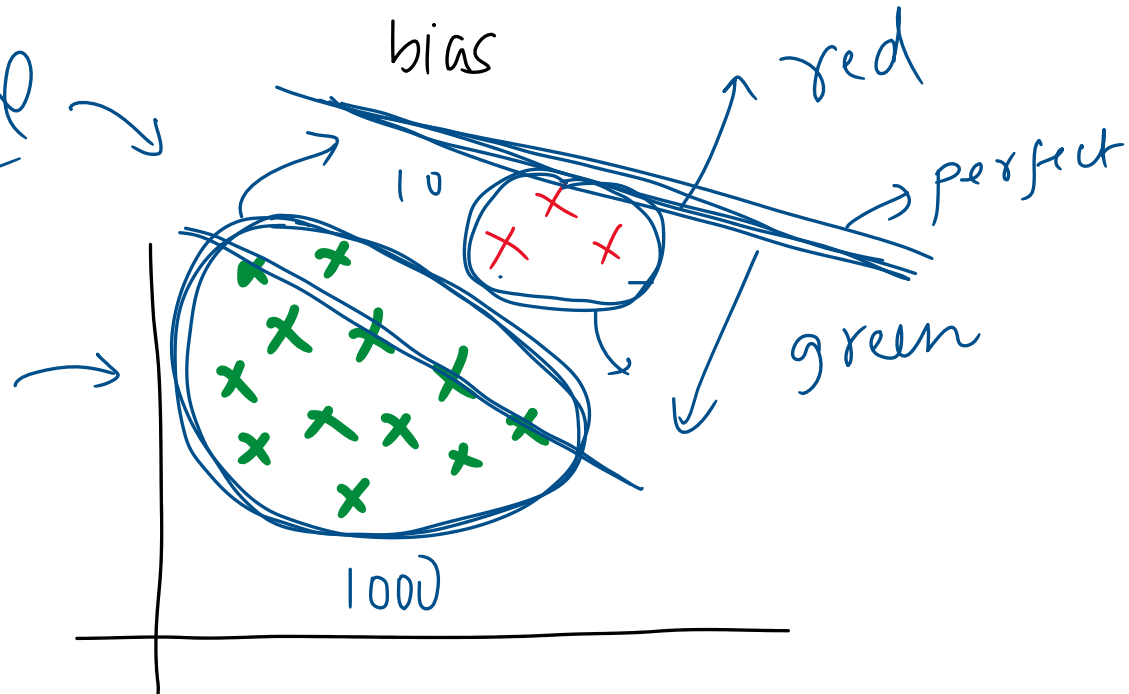
f1 score = harmonic mean (recall, precision)

24 April 2024 08:35

read ]  
pre ]

ml  $\rightarrow$  data  $\rightarrow$  [majority]  
 $\searrow$  minority  $\rightarrow$

ml  
internal



990  
-----  
100%

99 accuracy

# [Why Imbalanced Data is a huge deal?] → imp

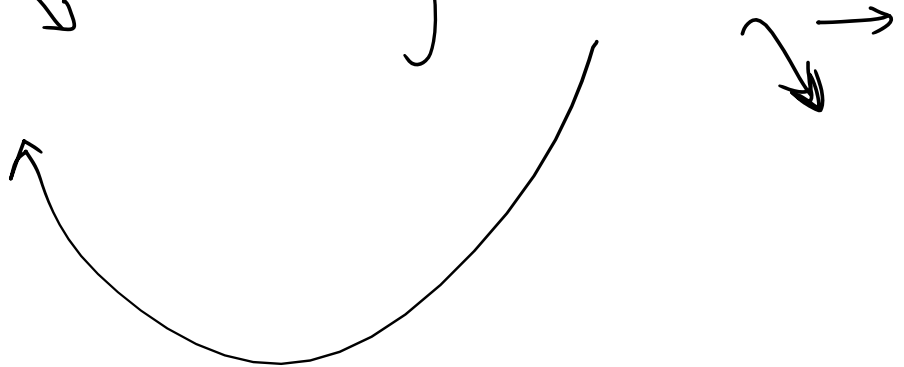
24 April 2024 08:35

Where can we find such data

- 1. Medical and healthcare
- 2. Fraud Detection
- 3. Manufacturing and Quality control
- 4. Customer Churn
- 5. Cybersecurity

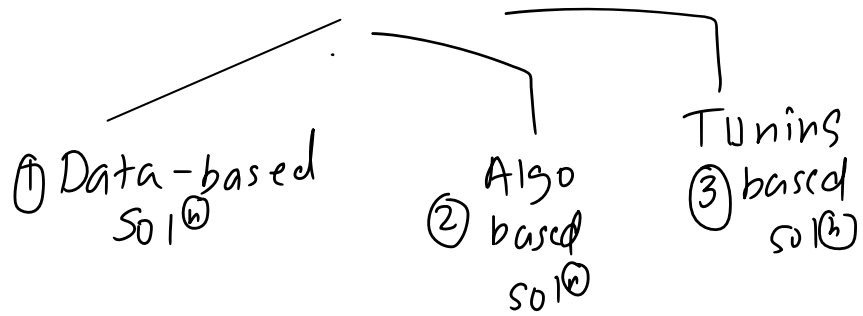
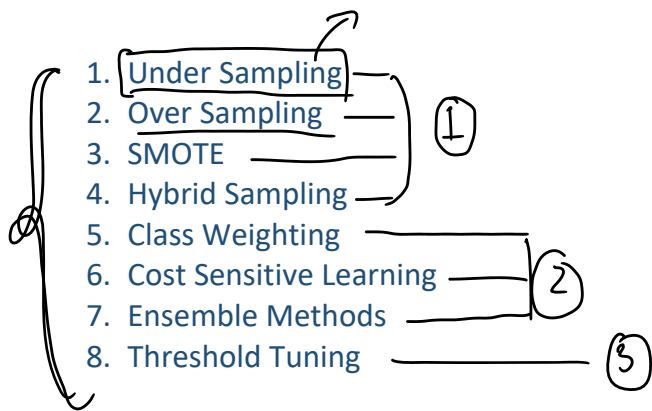
disease → imbalanced

ml models



# Techniques

24 April 2024 08:36



16:33

## Concept of Under sampling

## Why Is Under sampling Important?

- **Reduction of Bias:** Under sampling reduces the size of the majority class, thus reducing model bias and helping it learn from the minority class.
- **Faster Training:** With under sampling, the dataset becomes smaller, leading to faster training times.
- **Balanced Learning:** Models can learn from balanced data, allowing them to properly understand both classes.

samples

- **Loss of Information:** When you remove samples from the majority class, you might lose useful information, affecting model performance.
- **Overfitting:** A smaller dataset can lead to overfitting, especially if the minority class has limited examples.
- **Sampling Bias:** If under sampling is not done randomly, it can create a sampling bias, leading to a biased model.

There are several techniques for under sampling to balance imbalanced datasets:

- **Random Under sampling:** This technique involves randomly removing examples from the majority class to create balance.
- **Cluster-based Under sampling:** In this approach, examples from the majority class are grouped into clusters, and then some examples are chosen from each cluster to be removed.
- **Tomek Links:** This technique removes examples from the majority class that are close to the minority class, helping to clear the class boundaries.

balance

## Undersampling

undersampling

balanced

Original data

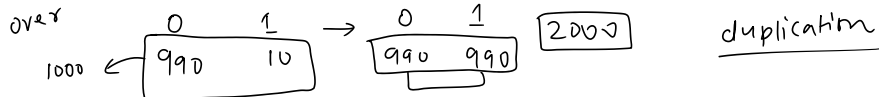
1





## 2. Over Sampling

24 April 2024 15:54



Oversampling is a technique used to balance imbalanced datasets. When a dataset has unequal class distribution, meaning one class (the majority class) has more samples and another class (the minority class) has fewer samples, oversampling is used to increase the size of the minority class to create a balance.

### Concept of Oversampling

The concept of oversampling is to increase the number of examples in the minority class to bring it closer to the majority class in terms of size. The objective is to provide balanced training data to machine learning algorithms so that they can learn patterns from both classes.

### Why Is Oversampling Important?

In imbalanced datasets, models often favour the majority class because it has more samples. Oversampling increases the size of the minority class, allowing models to learn from it properly. Here are some important reasons why oversampling is popular:

- **Improved Learning:** Oversampling provides more examples from the minority class, allowing models to learn more effectively.
- **Reduced Bias:** Oversampling creates a balanced dataset, reducing the model's bias toward the majority class.
- **Enhanced Performance:** A balanced dataset often improves model performance because it represents both classes equally.

### Risks of Oversampling

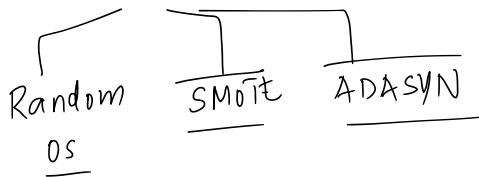
However, there are risks associated with oversampling:

- **Overfitting:** When minority class examples are oversampled, there is a risk of overfitting because the model may learn too much from the specifics of these examples.
- **Increased Data Size:** Oversampling increases the dataset's size, requiring more computational resources and potentially slowing down training.
- **Redundancy:** Oversampling often involves duplicating some examples, leading to redundancy in the dataset.

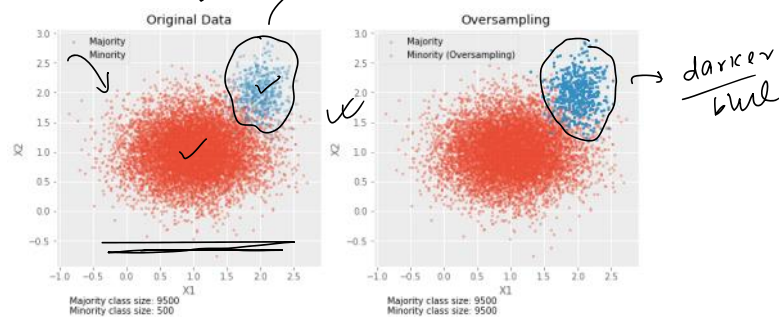
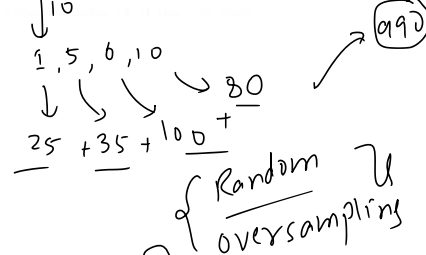
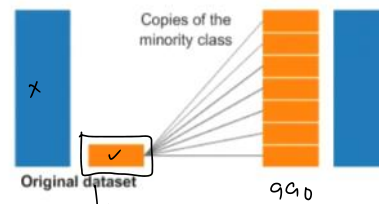
### Oversampling Techniques

To balance imbalanced datasets, there are several oversampling techniques:

- **Random Oversampling:** This technique involves randomly duplicating examples from the minority class to match the size of the majority class.
- **SMOTE (Synthetic Minority Over-sampling Technique):** This is an advanced technique where new synthetic examples are created based on existing minority class examples. The goal is to avoid duplication and bring diversity to the dataset.
- **ADASYN (Adaptive Synthetic Sampling):** This is a variation of SMOTE that creates more synthetic examples where the model has more difficulty.



### Oversampling



### 3. SMOTE

24 April 2024 16:38

oversampling minority  $\uparrow \rightarrow$  ROS

SMOTE stands for "Synthetic Minority Over-sampling Technique." It's a technique used to increase the size of the minority class in imbalanced datasets. The goal of SMOTE is to increase the number of examples in the minority class by creating synthetic data points, providing more balanced data for machine learning models.

#### Concept of SMOTE

The main concept of SMOTE is to create synthetic examples among existing examples in the minority class. This means generating new samples artificially to balance the data, providing sufficient training data for the model to learn the patterns in the minority class effectively.

#### SMOTE operates in the following way:

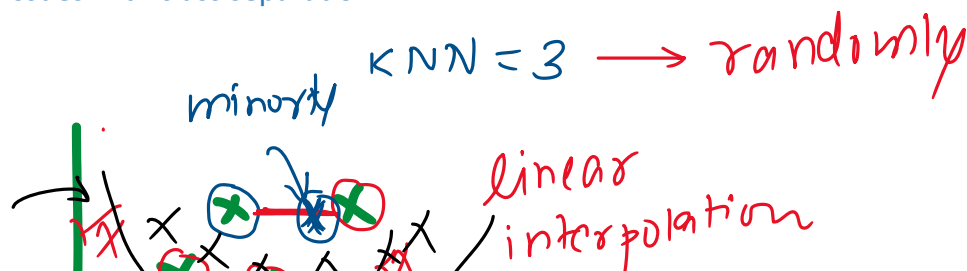
- **Using Nearest Neighbours:** SMOTE begins by finding the nearest neighbours for each example in the minority class. These nearest neighbours belong to the same class.
- **Creating Synthetic Examples:** For each example in the minority class, a few nearest neighbours are randomly selected. A synthetic example is then generated between the original example and the selected neighbour. Linear interpolation is used to create the synthetic example.
- **Creating a Balanced Dataset:** In this way, new synthetic examples are created for the minority class, bringing its size closer to or matching that of the majority class.

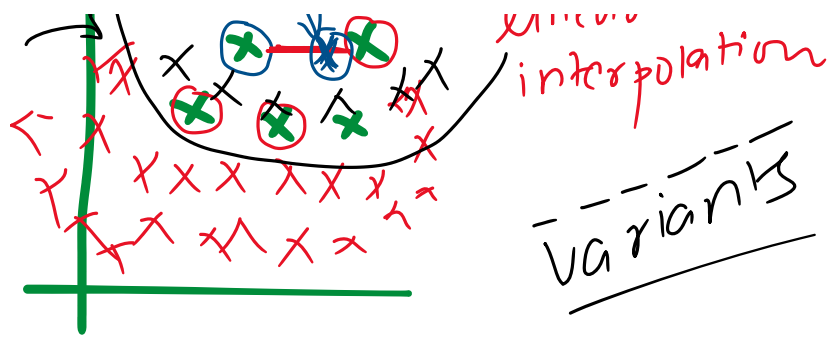
#### Benefits of SMOTE

- **Reduces Overfitting:** Because SMOTE creates synthetic examples, it avoids duplication, reducing the risk of overfitting.
- **Improves Learning:** The increase in the size of the minority class due to synthetic examples provides more training data, enabling the model to better learn the patterns in the minority class.
- **Enhances Model Performance:** Balanced datasets generally lead to better model performance, as they represent both classes more accurately.

#### Risks of SMOTE

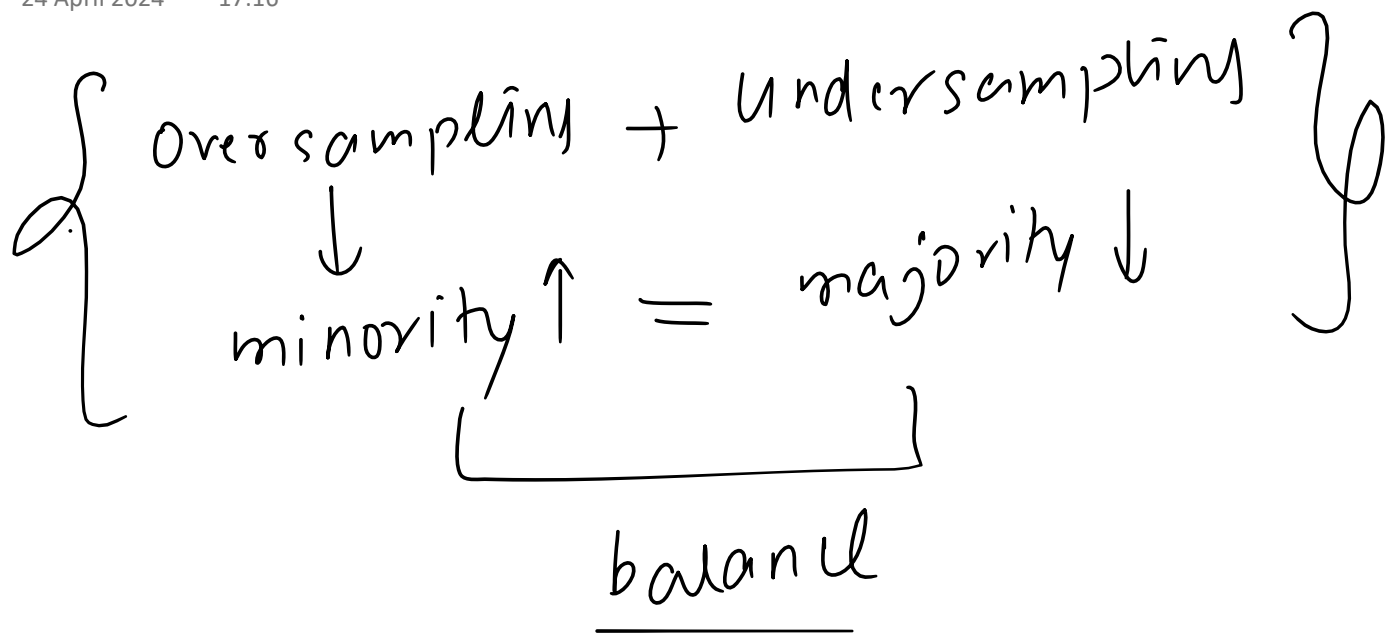
- **Boundary Issues:** Sometimes, SMOTE creates synthetic examples near decision boundaries, which can skew the boundary.
- **Complexity:** SMOTE's process is somewhat complex and can have high computational overhead.
- **Class Separation:** In some cases, over-creation of synthetic examples in the minority class can cause issues with class separation.





## 4. Hybrid Sampling

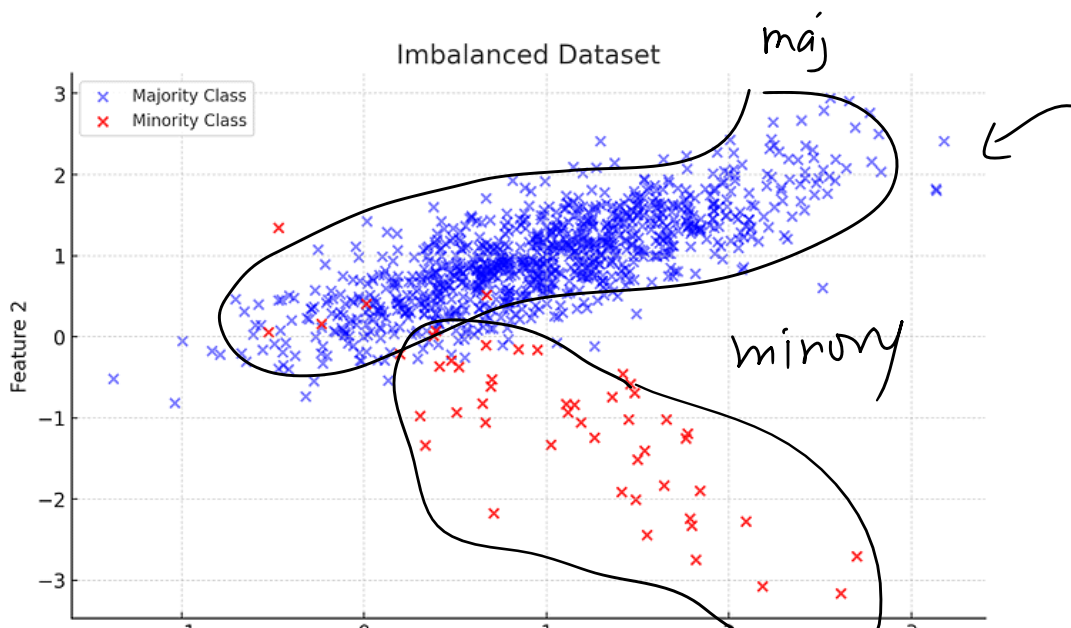
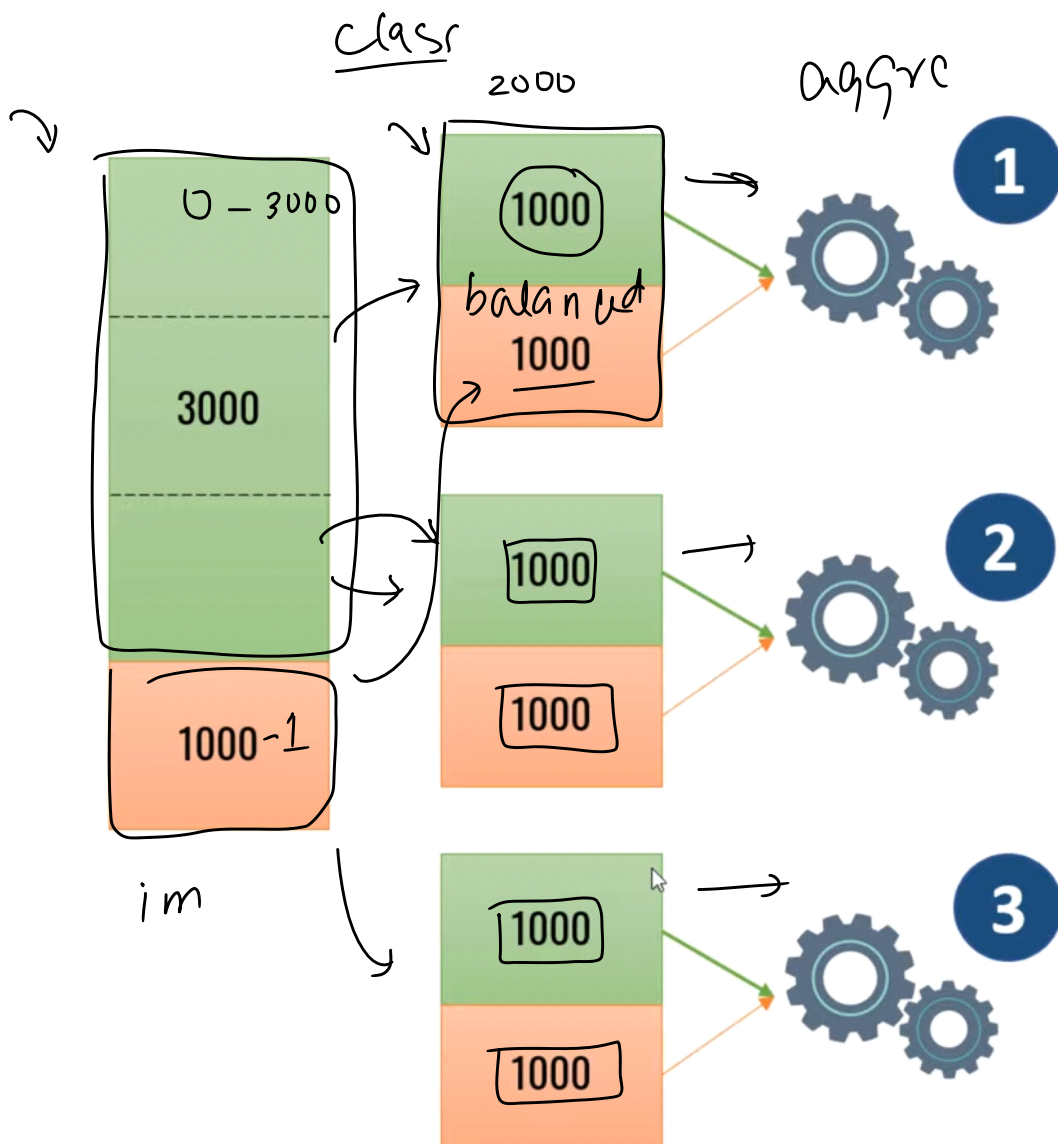
24 April 2024 17:16

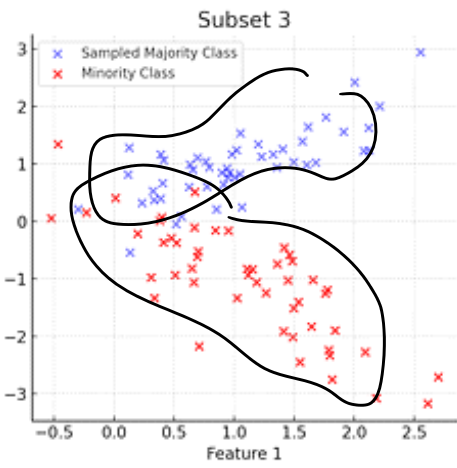
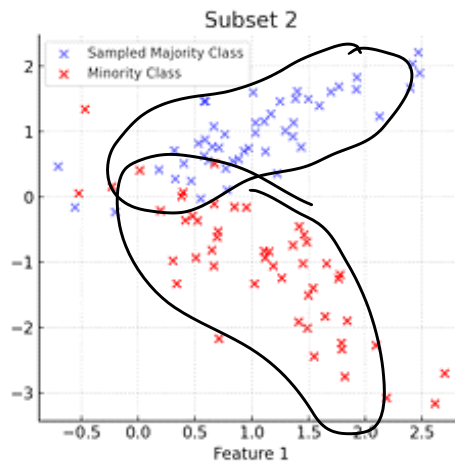
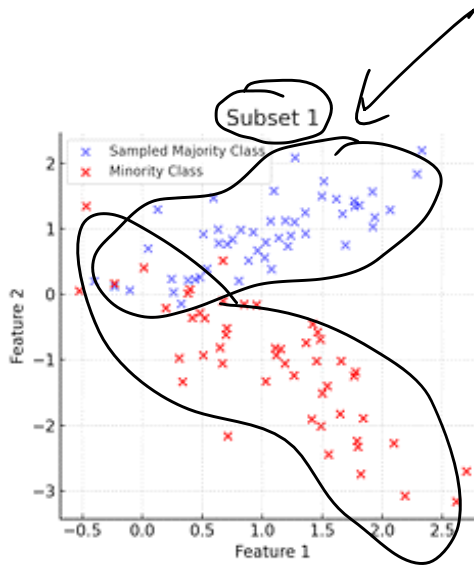
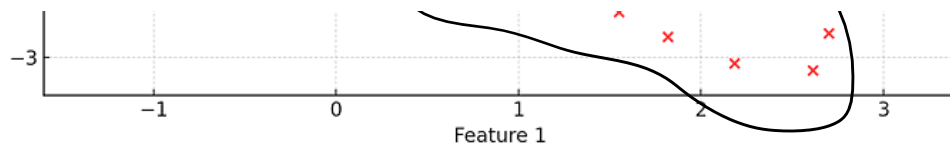


## 5. Ensemble Methods

24 April 2024 17:48

ensemble RF → Balanced RF





↓  
dt

↓  
dt

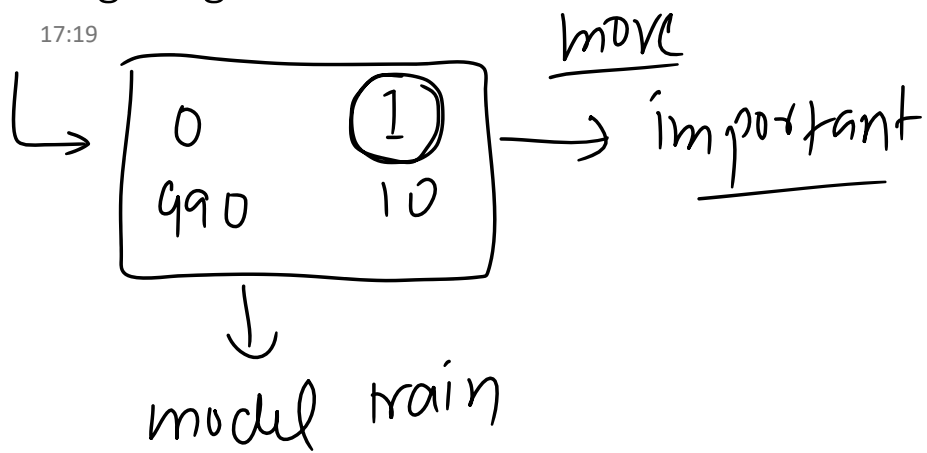
↓  
dt

Balanced  
Rf

ensemble

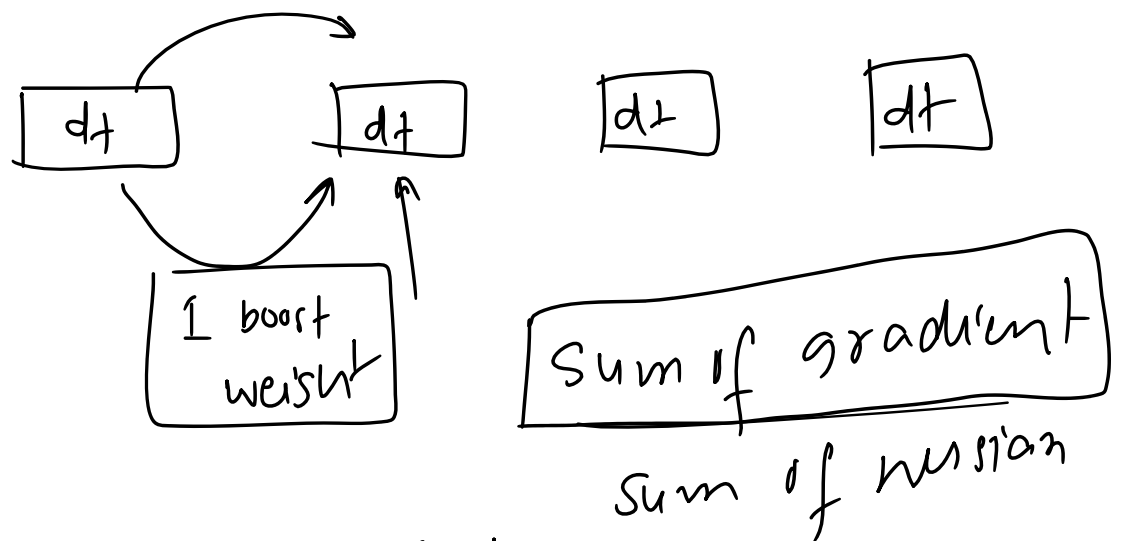
## 6. Class Weighting →

24 April 2024 17:19



$$\text{class weight} = \left\{ \frac{0:1}{1}, \frac{1:20}{} \right\}$$

## Boosting



$$\textcircled{1} \rightarrow \text{gradient} \times \textcircled{w}$$

$$\textcircled{0} \rightarrow \text{gradient} \times \textcircled{1}$$

$$\textcircled{1}$$

$$\text{hessian} \times w$$

$$\textcircled{0}$$

$$\text{hessian} \times 1$$

## 7. Cost Sensitive Learning

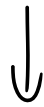
24 April 2024 17:48

model     $1 \uparrow$      $0 \downarrow$



diff / modified

loss function



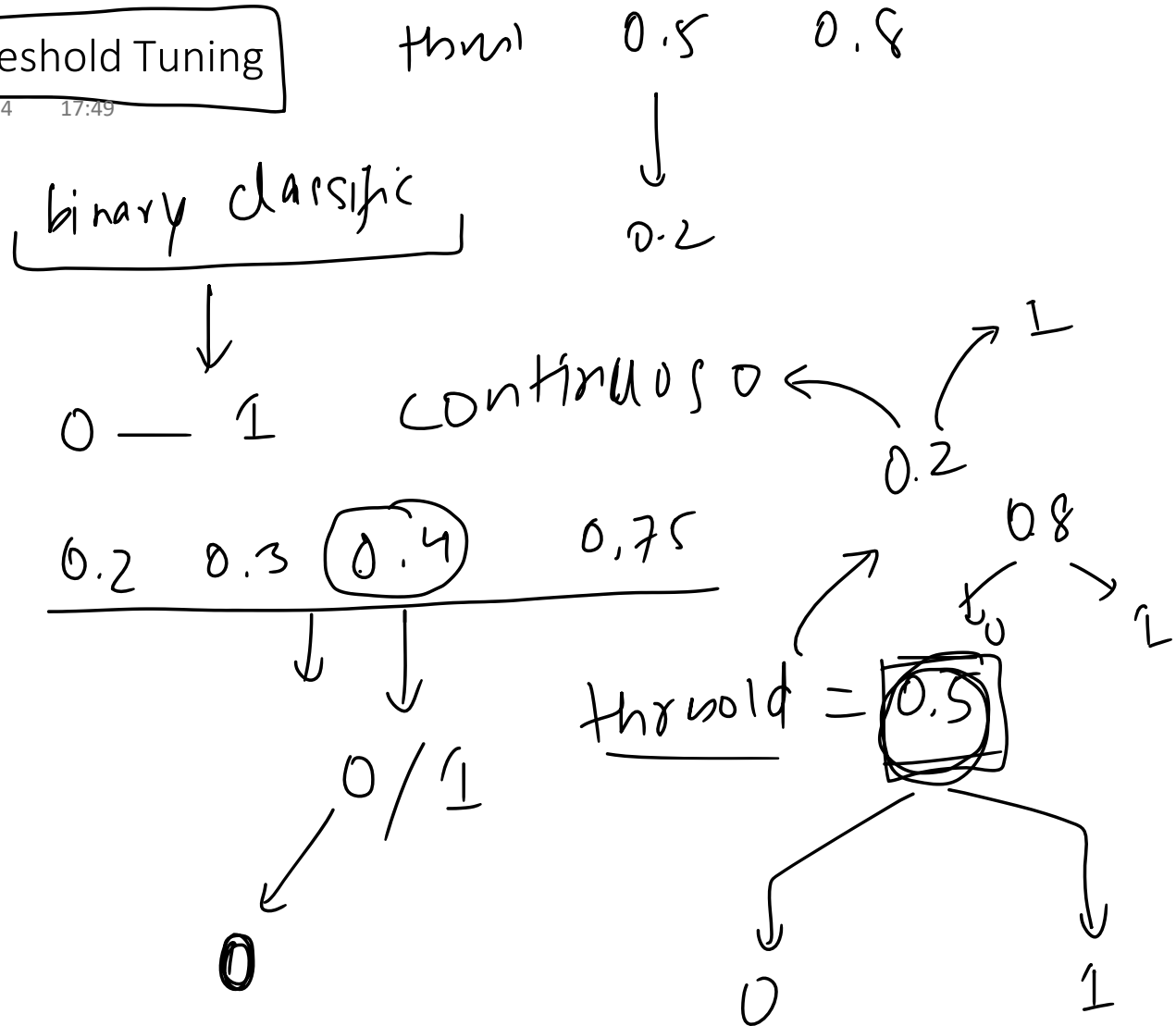
error  $\rightarrow$  ① (high)

error  $\rightarrow$  0 (low)



# 8. Threshold Tuning

24 April 2024 17:49



# List of All Techniques

24 April 2024 15:33