# Imbalanced Data

**Nguyen Khoa**

# CONTENT

## Introduction

Advantages:
- Can be applied to any existing learning tool
- The chosen models are biased to the goals of the user (because the data distribution was previously changed to match these goals), and thus it is expected that the models are more interpretable in terms of these goals.
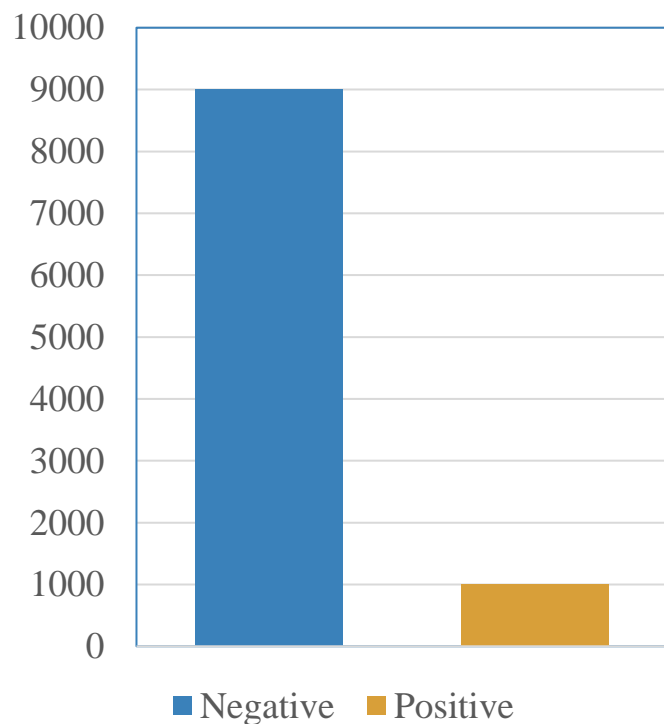
Disadvantages:
- Mapping the given data distribution into an optimal new distribution according to the user goals is not easy.
- It was proved for classification tasks that a perfectly balanced distribution does not always provide optimal results

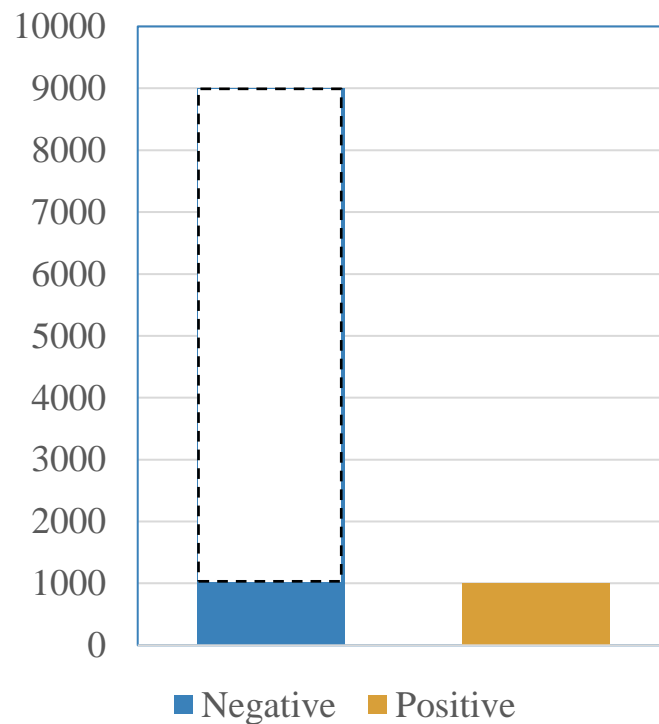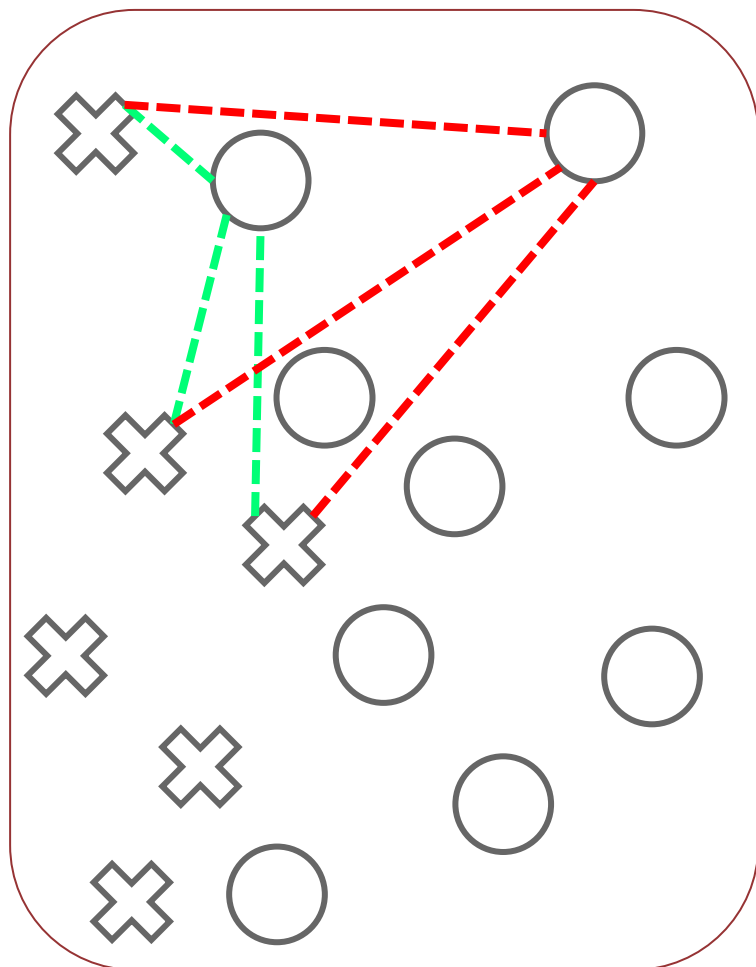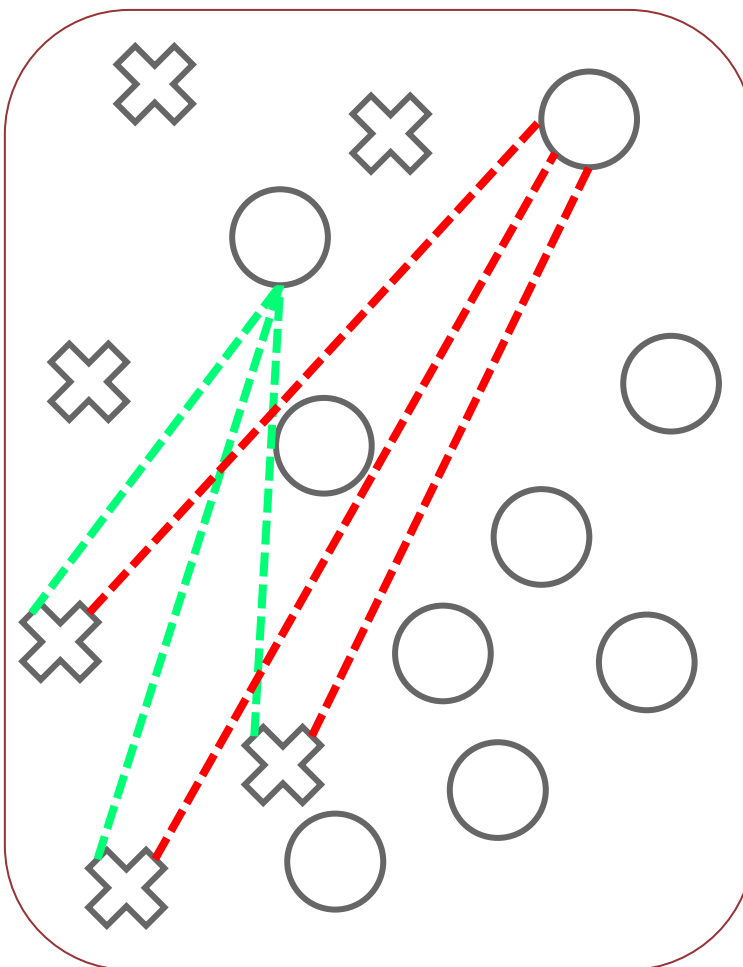**Re-Sampling / Distance-Based / NearMiss**



NearMiss1     NearMiss2     NearMiss1

**!** **Re-Sampling / Data Cleaning / Tomek Links**



**Tomek Links**  **Approach 1: Remove All**  **Approach 1: Remove Majority**

**Re-Sampling / Synthesising New Data / SMOTE**



**Majority Class Samples**

**Minority Class Samples**

**Randomly Selected Minority Class Sample ($x_i$)**

**n K-nearest Neighbors of $x_i$**

**Randomly Selected Sample from n K-nearest Neighbors**

**Generated Instance**

**Introduction**

Advantages:
- The user goals are incorporated directly into the models
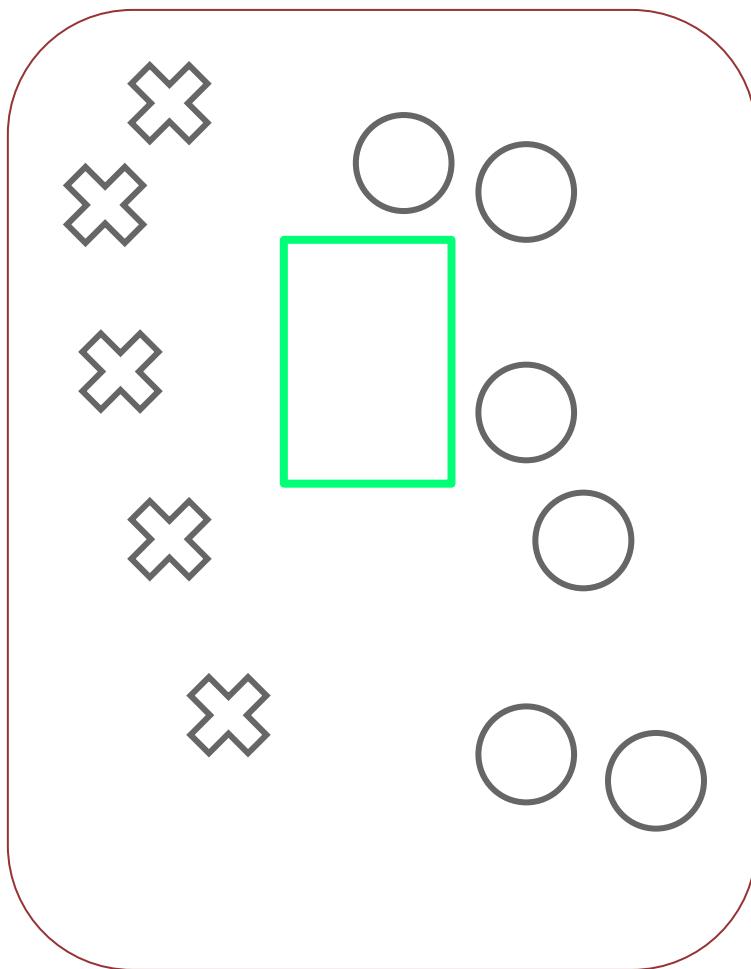
Disadvantages:
- restricted in his choice to the learning algorithms that have been modified to be able to optimise his goals, or has to develop new algorithms for the task
- may be necessary to introduce further modifications in the algorithm which may not be straightforward
- requires a deep knowledge of the learning algorithms implementations

**Transfer Learning**

**!** **Loss / MFE & MSFE**

$$FPE = \frac{1}{N}\sum_{i=1}^{N}\sum_{n}\frac{1}{2}\left(d_n^{(i)} - y_n^{(i)}\right)^2$$
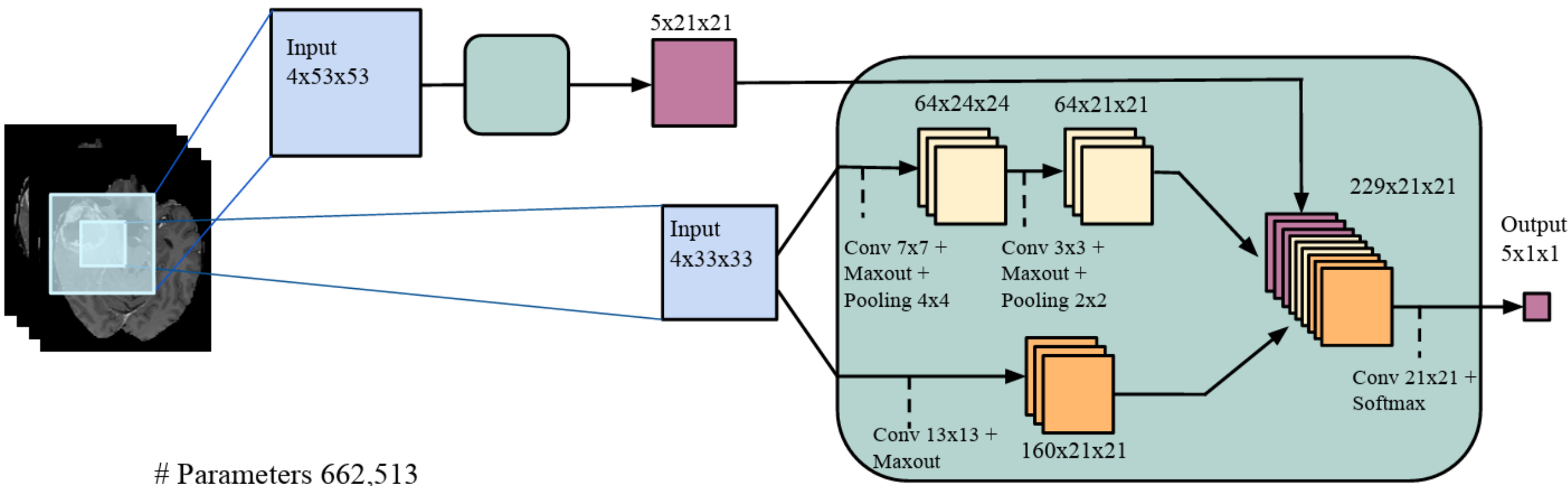
$$FNE = \frac{1}{P}\sum_{i=1}^{N}\sum_{n}\frac{1}{2}\left(d_n^{(i)} - y_n^{(i)}\right)^2$$

$$MFE = FPE + FNE$$

$$MSFE = FPE^2 + FNE^2$$

$MFE$: Mean False Error

$MSFE$: Mean Square False Error

$FPE$: Mean False Positive Error

$FNE$: Mean False Negative Error

$N$: the numbers of samples in negative class

$P$: the numbers of samples in positive class

$d_n^{(i)}$: the desired value of $i^{th}$ sample on $n^{th}$ neuron

$y_n^{(i)}$: the corresponding predicted value of $d_n^{(i)}$

11

**Optimizer / ImbSAM**



Figure 1: The visualization of separate loss landscape for *head* and *tail* classes in class-imbalanced recognition, optimized by SGD [6], SAM [17] and our ImbSAM respectively.

## Introduction

Advantages:
- Not necessary to be aware of the user preference biases at learning time
- Be applied to different deployment scenarios without the need of re-learning the models
- Any standard learning tool can be used

Disadvantages:
- The models do not reflect the user preferences
- The models interpretability is meaningless

## Thresholding

```
[2]    1 # Let's assume you have a classifier that outputs the following probabilities for a given datapoint x
       2 # These are just example values for demonstration purposes
       3 classifier_output_A = 0.3  # Probability that x belongs to Class A
       4 classifier_output_B = 0.7  # Probability that x belongs to Class B
       5
       6 # Number of instances in each class
       7 instances_A = 10
       8 instances_B = 70
       9
      10 # Total number of instances
      11 total_instances = instances_A + instances_B
      12
      13 # Calculating the prior probabilities for each class
      14 prior_prob_A = instances_A / total_instances
      15 prior_prob_B = instances_B / total_instances
      16
      17 # Adjusting the classifier's probabilities by the prior probabilities
      18 adjusted_prob_A = classifier_output_A / prior_prob_A
      19 adjusted_prob_B = classifier_output_B / prior_prob_B
      20
      21 # Normalizing the adjusted probabilities so that they sum up to 1
      22 sum_adjusted_probs = adjusted_prob_A + adjusted_prob_B
      23 normalized_prob_A = adjusted_prob_A / sum_adjusted_probs
      24 normalized_prob_B = adjusted_prob_B / sum_adjusted_probs
      25
      26 (adjusted_prob_A, adjusted_prob_B), (normalized_prob_A, normalized_prob_B)

   ((2.4, 0.7999999999999999), (0.75, 0.25))
```

14

**Notion**

https://transparent-sesame-adf.notion.site/Imbalanced-4310cb6fbcfb4aa493145ef244c02f43

# Thanks!

## Any questions?