



# **TECHNIQUE TO HANDLE IMBALANCED DATASETS**

**Manohar Kuse**

**Kumar Harsh Srivastava**

**Abhishek Dubey**

**Aniroop Mathur**

# PROBLEM DESCRIPTION

Class A



Class B



$$| N_A | \gg | N_B |$$



# MOTIVATION

- Rare Patterns give skewed datasets
  - Cancer Cell Classification
  - Gene Profiling
  - Credit Card Frauds
- Minority class considered as noise
- Improper training of classifier due to skew
- Required to balance the dataset by artificial samples

# IDEA

- Artificial instances by combining k-minority samples at a time
- $F \rightarrow$  Feature vectors of minority class
- $M \rightarrow$  Number of samples in minority class
- $R = \{r \mid 1 < r < M, r \text{ is a random number}\}$   
such that,  $|R| = k$

$$s_{new,p} = 1/k \sum_{x_i \in R} F(x_i)$$



## PSEUDO CODE

- Inputs : F, k, NE (number of artificial samples)
- Output : E (Artificial samples)

```
For i = 1 to NE
  <tmp> = < 0 >
  For j = 1 to k
    r = random( 1, M )
    <tmp> = <tmp> + <F(r)>
  End
  E(i) = <tmp> / k
End
```



## REFERENCES

- **Handling imbalanced datasets: A review.** S. Kotsiantis, D. Kanellopoulos, P. Pintelas.
- **SMOTE: Synthetic Minority Over-sampling Technique.** N. Chawla, K. Bowyer, L. Hall, W. Kegelmeyer.
- **Off-line, Handwritten Numeral Recognition by Perturbation Method.** T. Ha, H. Bunke

