



# SMOTE VARIANT FOR UNBALANCED DATA IN CLASSIFICATION PROBLEMS

Valeria De Stasio, Christian Faccio, Andrea Suklan, Agnese Valentini

January 26, 2025

- Introduction to the problem
- Data generation / Assessment method
- SMOTE-DIRICHLET
- Results

# Data Generation Process

Parameter	Values	Description
$n_{\text{train}}$	600, 1000, 5000	Train set size
$n_{\text{test}}$	600	Test set size
$\pi$	0.10, 0.05, 0.025	Proportion of rare examples
$IR$	9, 19, 39	Imbalance ratio

Table: Parameters of the simulation

$$(\mathbf{X}, y) \text{ s.t. } \begin{cases} \mathbf{x} \sim N_2 \left( \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right) & \text{if } y = 0, \\ \mathbf{x} \sim N_2 \left( \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 & -0.5 \\ -0.5 & 1 \end{pmatrix} \right) & \text{if } y = 1. \end{cases}$$

# Data Generation

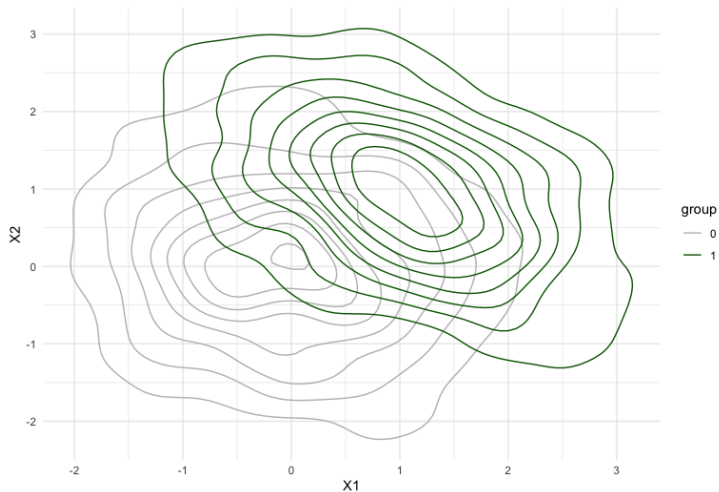


Figure: Contour plot of the data distribution

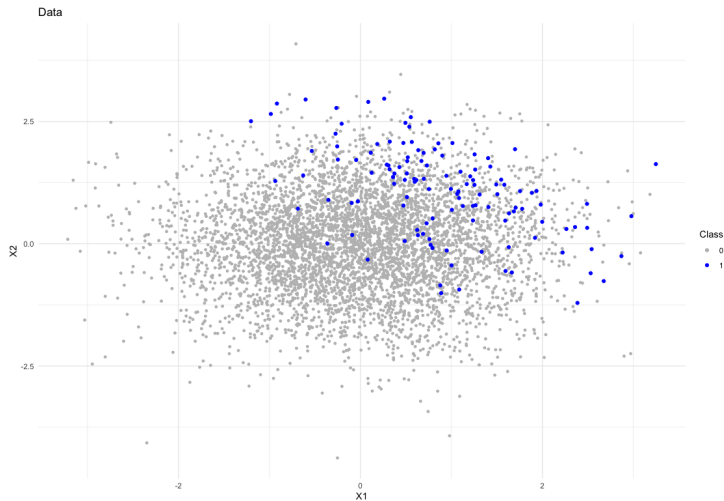


Figure: Data generated with  $size = 5000$  and  $prop = 0.025$

# Assessment Method

- 100 Static Train/Test Division (test size = 600)
- Metrics: Balanced Accuracy, F1 Score
- Models: Decision Tree, Logistic Regression
- Comparisons with Boxplots

# SMOTE

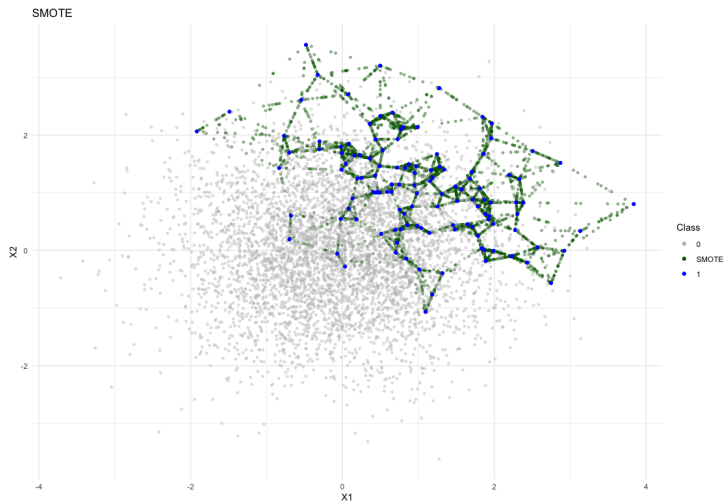
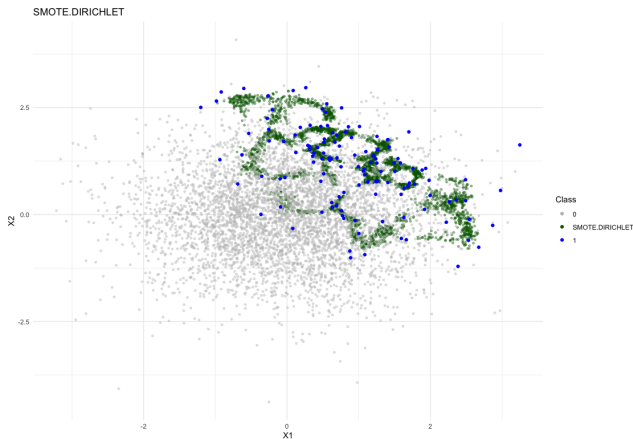


Figure: Synthetic data generated with SMOTE

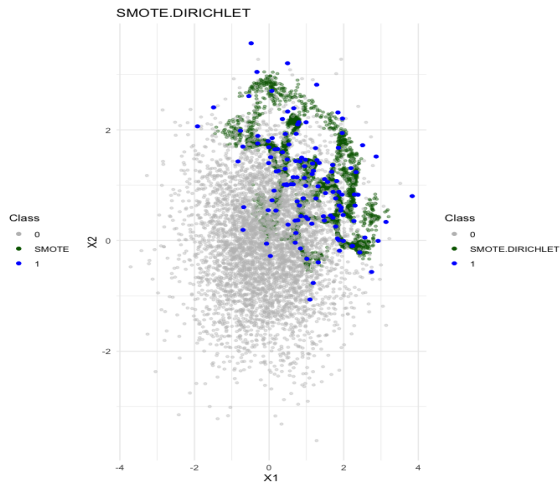
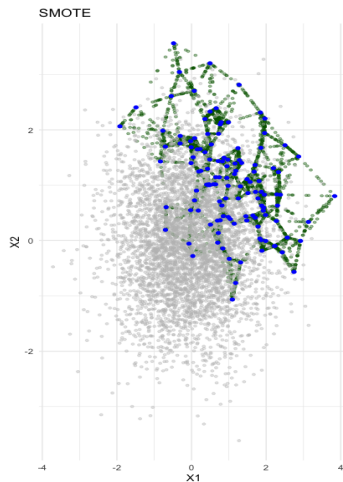
# Dirichlet SMOTE



**Figure:** Synthetic data generated with Dirichlet SMOTE using Dirichlet distribution:  $P(p|\alpha) \sim \text{Dir}(\alpha_1, \dots, \alpha_k) = \frac{\Gamma(\sum_j \alpha_j)}{\prod_j \Gamma(\alpha_j)} \prod_{j=1}^k p_j^{\alpha_j-1}$  [3].  $\alpha = 1$ ,  $k = 3, 5$ .



# Comparison between SMOTE and Dirichlet SMOTE



# References



Chawla, Nitesh V et al. (2002). “SMOTE: synthetic minority over-sampling technique”. In: *Journal of artificial intelligence research* 16, pp. 321–357.



Lunardon, Nicola, Giovanna Menardi, and Nicola Torelli (2014). “ROSE: a package for binary imbalanced learning”. In: *R journal* 6.1.



Matharaarachchi, Surani, Mike Domaratzki, and Saman Muthukumarana (2024). “Enhancing SMOTE for imbalanced data with abnormal minority instances”. In: *Machine Learning with Applications* 18, p. 100597.



Menardi, Giovanna and Nicola Torelli (2014). “Training and assessing classification rules with imbalanced data”. In: *Data mining and knowledge discovery* 28, pp. 92–122.