



UNIVERSITÀ  
DEGLI STUDI  
DI TRIESTE

---

# SMOTE VARIANT FOR UNBALANCED DATA IN CLASSIFICATION PROBLEMS

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

January 25, 2025

# Content

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

# Data Generation Process

Parameter	Values	Description
$n_{\text{train}}$	250, 1000, 5000	Train set size
$n_{\text{test}}$	250	Test set size
$\pi$	0.1, 0.05, 0.025	Proportion of rare examples

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

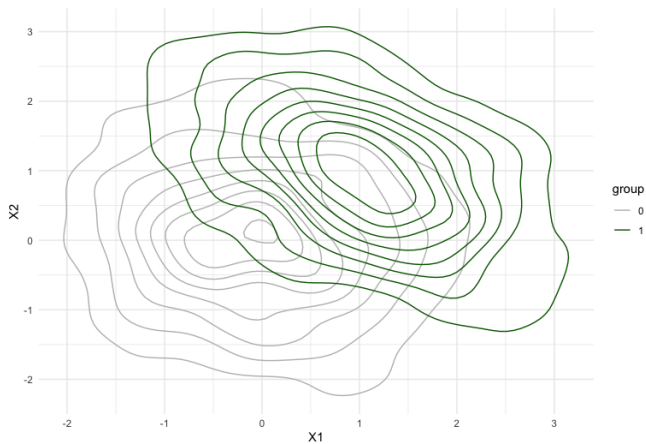






Figure: Contour plot of the data distribution

# 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.