

# synthetic minority oversampling using multivariate gaussian distribution

namagiri December 2023

## 1 Abstract

Highly imbalanced data is a common phenomenon from cases of fraudulent transactions to several medical datasets. Without rebalancing the dataset if it fits into any learning algorithm will result in a highly biased model towards the majority class to address this issue we have employed a technique of oversampling using a new smote(synthetic minority oversampling ) strategy where  $k$  nearest neighbours are picked from the minority class and fitted to multivariate gaussian distribution to generate points which populate regions with high probability

## 2 Methodology

Partially guided synthetic minority oversampling using Gaussian kernel interpolation.

Let's consider the matrix  $X$ :

$$X = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{a1} & a_{a2} & \cdots & a_{na} \end{bmatrix}$$

The dimensions of matrix  $X$  are  $n \times a$ , where  $n$  refers to rows/instances and  $a$  refers to attributes.

Let  $t$  be the target attribute and  $X[t] = \{0, 1\}$ , representing a binary classification task.

$X = \text{Maj} \cup \text{Min}$ , where Maj is the majority class and Min is the minority class.

$\text{Maj}[t] = \{0\}$  and  $\text{Min}[t] = \{1\}$ .

Let  $\theta = \frac{n(\text{Min})}{n(\text{Maj})}$ , where  $\theta \in [0, 1]$ .

$n()$  denotes the cardinality of the class, and  $\theta \ll 1$ .

---

**Algorithm 1** Selecting ensemble classifier

---

- 1: **Input:** Minority data set  $Min$ , Majority data set  $Maj$ , ratio balancer  $\delta$
  - 2: **Output:** Ensemble classifier  $E$
  - 3: bounding the number of subsets  $b = \lceil \frac{n(Maj)}{\delta * n(Min)} \rceil$
  - 4: Selecting  $b$  subsets from  $Maj$  |  $\bigcup_{i=1}^b Maj(i) = Maj$  and  $\bigcap_{i=1}^b Maj(i) = \emptyset, \forall i, n(Maj(i)) > \delta * n(Min)$
  - 5: Number of synthetic samples to be generated
  - 6: sample number =  $\min_{i \in [1, b]} n(Maj(i)) - n(Min)$
  - 7:  $Gen = Algorithm_2(Min, Samplenumber)$
  - 8:  $Min = Min \cup GEN$
  - 9:  $C_i(Maj(i), Min)$
  - 10: Getting  $b$  classifiers and trained to evaluate the best classifier  $BEST = decide(\{C_1, C_2, C_3, \dots, C_b\})$
  - 11: **Return:**  $BEST$
- 

---

**Algorithm 2** Generating synthetic samples

---

- 1: **Input:** Minority class  $Min$ , Sample number
  - 2: **Output:** Synthesized samples  $GEN$
  - 3: Start  $s = sample\ size$  and  $GEN \leftarrow \phi$
  - 4: **for**  $i \in Min$  **do**
  - 5:    $l = \text{Sort}(k \text{ nearest neighbours of } i)$
  - 6:    $l\_first \leftarrow \text{First } \frac{l}{4} \text{ elements of data}$
  - 7:    $l\_last \leftarrow \text{Last } \frac{l}{4} \text{ elements of data}$
  - 8:    $l_1 = l\_last \cup l\_first$
  - 9:    $f \leftarrow \text{Multivariate gaussian distribution}(l_1)$     $\triangleright$  Interpolated surface
  - 10:   **Evaluation:**
  - 11:    $l\_middle \leftarrow \text{Middle } \frac{l}{2} \text{ elements of data}$
  - 12:   **for**  $x$  in  $l\_middle$  **do**
  - 13:      $y \leftarrow \text{Evaluate}(f, x)$     $\triangleright$  Evaluate the function at the middle elements
  - 14:   **end for**
  - 15:    $s \leftarrow s - n(y)$
  - 16:    $GEN \leftarrow GEN \cup y$
  - 17:    $Min \leftarrow Min - l$
  - 18: **end for**
  - 19: **Return:**  $GEN$
-

## fitting the l1 vector into the distribution

Let  $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$  be a vector-valued random variable with a multivariate normal distribution. It has mean  $\boldsymbol{\mu} \in R^n$  and covariance matrix  $\boldsymbol{\Sigma} \in S_n^{++}$ .

The probability density function (pdf) of  $\mathbf{X}$  is given by:

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{n/2} |\boldsymbol{\Sigma}|^{1/2}} \exp \left( -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right)$$

To generate samples from this distribution, we can use the Cholesky decomposition. The Cholesky decomposition expresses the positive definite matrix  $\boldsymbol{\Sigma}$  as the product of a lower triangular matrix  $\mathbf{L}$  and its transpose:

$$\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$$

The algorithm for generating samples is as follows:

1. Ensure that  $\boldsymbol{\Sigma}$  is positive definite.
2. Compute the Cholesky decomposition:  $\mathbf{L} = \text{cholesky}(\boldsymbol{\Sigma})$ .
3. Generate standard normal samples  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ .
4. Transform to multivariate normal samples:  $\mathbf{X} = \boldsymbol{\mu} + \mathbf{z}\mathbf{L}^T$ .