

# Outlier SMOTE Algorithm

Research

Name: Meenavalli Sri Sai Nitin Ray Vansi

BITID: 2020CT04508

Group: 12

## Objective:-

Class imbalance is a common issue in every Data set. Due to this data becomes biased to specific classes, training classifiers is becoming difficult, hence effecting its performance.

⇒ This issue is tackled by using over (or) under Sampling Techniques, the one which stood above all is the SMOTE algorithm, which generates synthetic samples of each minority class by oversampling of ~~minority~~ data <sup>each</sup> point by combination of existing minority class samples neighbours.

⇒ This may lead to overlapping of minority data samples. To overcome this issue authors come up with a modified version of SMOTE known as Outlier SMOTE.

⇒ Where in each data point is oversampled with respect to its distance from the other data points.

⇒ The data point which is farthest from the data point is given high priority and is oversampled more than its Counterparts.

## Data Pre-Processing:-

⇒ Cleaning the dataset by removing features with >90% null values, variables with ~~some~~ variance, and Counting Columns by taking row wise sum to identify the minority samples, and pass them to OUTLIER SMOTE.

⇒ This will return an euclidean distance matrix of new size, for n samples and m features.

• Constructing Oversampling Matrix each Sample is assigned probability weight  $P(0 \text{ to } 1)$ .

→ This gives the sum of Columns gives the distance of samples from minority classes

⇒ Now the oversampling matrix may contain all non-integers, So normalizing it

→ Rounding off values of normalized matrix upto 2 decimal points.

Data Mining:-

→ Using Binary Logistic Regression Classification algorithm and gauging its performance by using different Sampling rates.

Metric:-

⇒ The main goal of outlier SMOTE algorithm is to reduce false positives. So we need Confusion Matrix, Precision, recall, F-measure

⇒ Precision = Correctly classified Positive / Total positive Samples.

⇒ Recall = No. of " " " " No. of Observations in positive class.

⇒ F-measure = Harmonic mean of Precision & recall.

Visualization:-

⇒ The Authors used SHAP algorithm which is a Game theoretic approach model.

⇒ Denoting which features are important or which features dominant

⇒ Prediction of an instance  $X$  by Computing the Contribution of each feature to the prediction

⇒ The X-axis contains the impact of features on result (-ve means less impact, +ve big)

⇒ The Y-axis has all features that have been observed.