

The detailed description of all selected datasets and the types of outliers

I: The detailed description of all selected datasets

The investigated datasets are all selected from publicly accessible repositories and are commonly used for the evaluation of outlier detection methods. For example, the “BreastW” dataset is fetched from the domain of medical diagnosis, where the malignant samples are treated as outliers to be identified. The detailed information about these datasets is described in Table 1.

TABLE 1: The detailed descriptions of all selected datasets

	Description
Arrhyth	Samples from the minority classes 3, 4, 5, 7, 8, 9, 14, and 15 are grouped into outliers, while the remaining samples are considered as inliers.
Autos	Samples from “-2” and “-1” classes are regarded as outliers, and samples from other classes are considered as inliers.
Breast	Samples in the “malignant” class are treated as outliers, while samples in the “benign” class are considered as inliers.
Cardio	The “pathologic” class is downsampled to 176 samples, and these samples are regarded as outliers. Samples in the “normal” class are considered as inliers, and samples in the “suspect” class are all removed.
Cardioto	The classes “2” and “3” are downsampled to 33 samples and selected as outliers, while other classes are all considered inliers.
Carpet	This is a category from the commonly used MVTec AD datasets and contains both normal and defective carpet samples.
Chess	The classes “nowin” is downsampled to 227 samples and treated as outliers, while other samples from the remaining class are considered as inliers.
Hepat	Samples in the “hepatitis” class are regarded as outliers, and samples in the “non-hepatitis” class are considered as inliers.
Iono	The classes “b” is downsampled to 24 anomalous samples, where samples from the remaining class are recognized as inliers.
Iris	The “iris-virginica” class is downsampled to 11 anomalous samples, and samples from other classes are considered as inliers.
Mammo	Samples in the “calcification” class are regarded as outliers, and other samples are considered as inliers.
Metal	This is a category from the commonly used MVTec AD datasets and contains both normal and defective metal samples.
Pen	Samples in the “0” class are regarded as outliers, and the remaining digit samples (1 to 9) are considered as inliers.
Pill	This is a category from the commonly used MVTec AD datasets and contains both normal and defective pill samples.
Sat	The “2” class is downsampled to 71 anomalous samples, while samples from other classes are combined to be inliers.
Spam	Samples in the “spam” class are regarded as outliers, and samples in the “non-spam” class are considered as inliers.
Thyroid	Samples in the “hyperfunction” class are treated as outliers, while samples in the normal and subnormal functions are considered as inliers.
WDBC	The “M” class is downsampled 39 samples and recognized as outliers, while other samples are considered as inliers.
Wine	The “1” class is downsampled to 10 samples and treated as outliers, while samples from the “2” and “3” classes are considered as inliers.
WPBC	Samples in the “R” (minority) class are recognized as outliers, and other samples are regarded as inliers.

II: The types of outliers in real datasets

Realistic datasets may entangle with multiple types of outliers, but decoupling these types of outliers requires expensive costs. To explore this problem, we visualized some datasets used in the experiments and selected some representative results. Figs. 1, 2, and 3 show global outliers, local outliers, and group outliers, respectively, in realistic datasets.

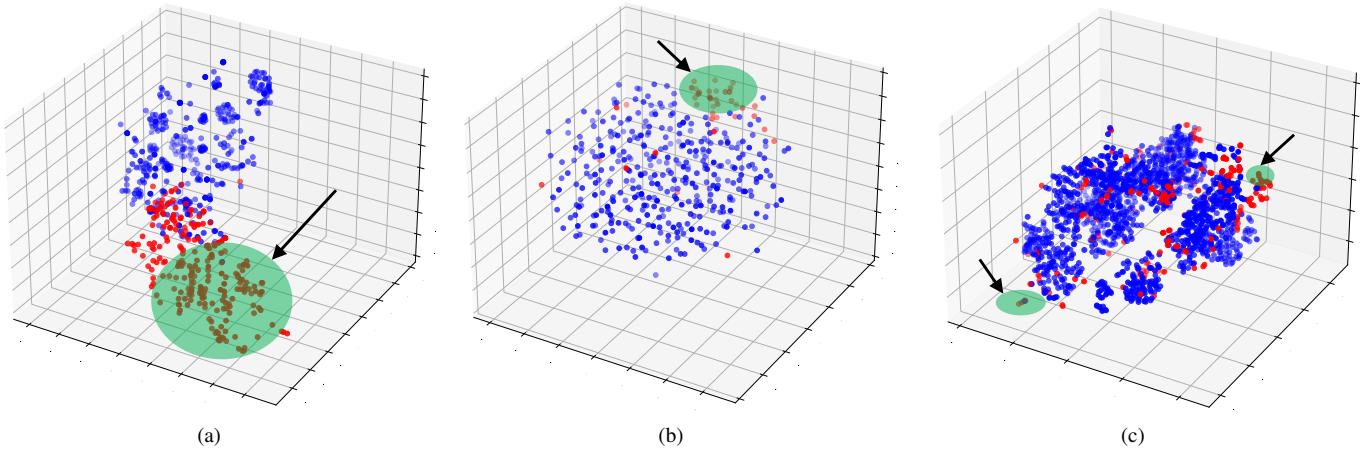


Fig. 1: The real global outliers visualized in 3-dimensional t -SNE space. a) Breast; b) WDBC; c) Chess.

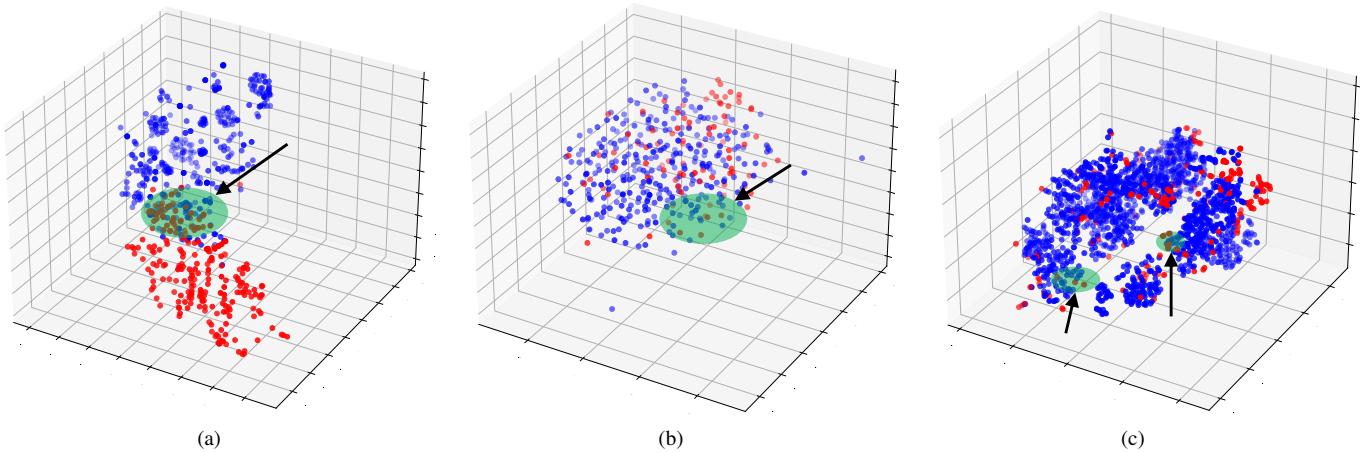


Fig. 2: The real local outliers visualized in 3-dimensional t -SNE space. a) Breast; b) Carpet; c) Chess.

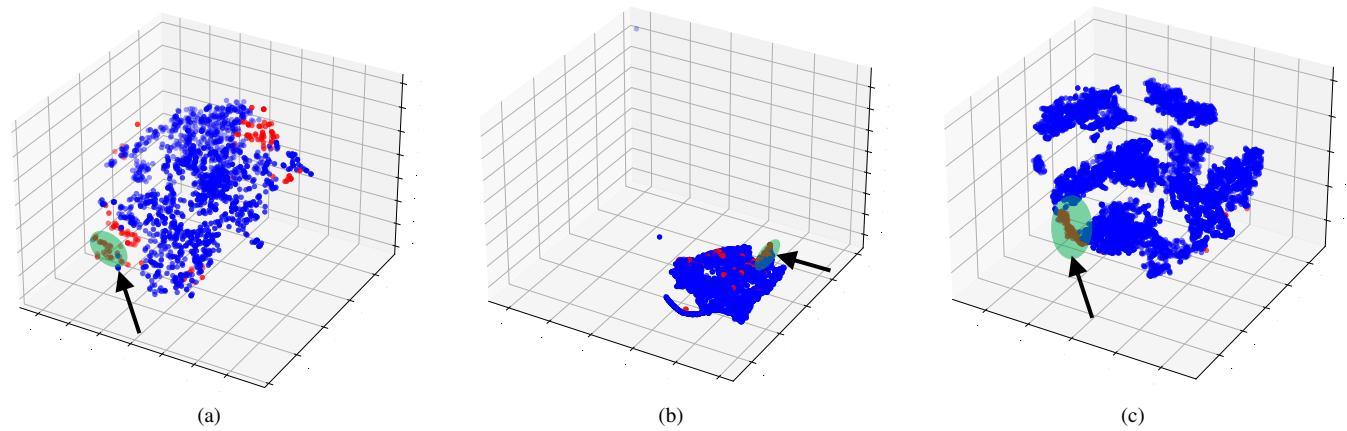


Fig. 3: The real group outliers visualized in 3-dimensional t -SNE space. a) Cardio; b) Mammo; c) Pen.

As shown in Figs. 1-3, global outliers, local outliers, and group outliers do exist in real datasets. Note that these types of outliers may coexist. For example, the “Breast” and “Chess” datasets contain both global and local outliers, while the “Cardio” and “Pen” datasets include global outliers and group outliers.