

3rd May Assignment

May 9, 2023

1 Assignment 85

Q1. What is the role of feature selection in anomaly detection?

Ans. The role of feature selection in anomaly detection is to identify and select the most relevant features that are useful in distinguishing between normal and anomalous data points. By reducing the dimensionality of the dataset, feature selection can improve the performance and efficiency of anomaly detection algorithms, reduce the risk of overfitting, and increase the interpretability of the results. Feature selection can be performed using various techniques such as correlation analysis, mutual information, principal component analysis (PCA), and recursive feature elimination (RFE).

Q2. What are some common evaluation metrics for anomaly detection algorithms and how are they computed?

Ans. Some common evaluation metrics for anomaly detection algorithms include precision, recall, F1-score, receiver operating characteristic (ROC) curve, and area under the curve (AUC). Precision measures the proportion of detected anomalies that are truly anomalies, while recall measures the proportion of true anomalies that are detected by the algorithm. F1-score is the harmonic mean of precision and recall, and provides a balanced measure of their performance. ROC curve plots the true positive rate (TPR) against the false positive rate (FPR) at different decision thresholds, while AUC measures the area under the ROC curve and provides an overall measure of the algorithm's ability to distinguish between normal and anomalous data points.

Q3. What is DBSCAN and how does it work for clustering?

Ans. DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a clustering algorithm that identifies clusters of data points based on their density in the feature space. DBSCAN works by defining a neighborhood around each data point and counting the number of points within the neighborhood. Data points that have a minimum number of neighbors within a specified radius are considered core points, and are used to form clusters. Data points that have fewer neighbors than the minimum threshold but are within the radius of a core point are considered border points and are assigned to the same cluster. Data points that have fewer neighbors than the minimum threshold and are not within the radius of any core point are considered noise points and are not assigned to any cluster. The key parameters of DBSCAN

are the radius of the neighborhood (eps) and the minimum number of neighbors (min_samples) required to form a core point.

Q4. How does the epsilon parameter affect the performance of DBSCAN in detecting anomalies?

Ans.The epsilon parameter in DBSCAN determines the radius of the neighborhood around each data point, and thus affects the ability of the algorithm to detect anomalies. A smaller value of epsilon leads to the formation of smaller clusters and more noise points, which may make it easier to detect anomalies that are isolated from the majority of the data.

Q5. What are the differences between the core, border, and noise points in DBSCAN, and how do they relate to anomaly detection?

Ans.In DBSCAN, core points are data points that have a minimum number of neighbors within a specified radius, while border points have fewer neighbors than the minimum threshold but are within the radius of a core point, and noise points have fewer neighbors than the minimum threshold and are not within the radius of any core point. Core points and border points can be considered as part of normal data, while noise points can be considered as anomalies.

Q6. How does DBSCAN detect anomalies and what are the key parameters involved in the process?

Ans.DBSCAN detects anomalies by labeling data points that are not assigned to any cluster as noise points, which can be considered as anomalies. The key parameters involved in the process are the radius of the neighborhood (epsilon) and the minimum number of neighbors (min_samples) required to form a core point.

Q7. What is the make_circles package in scikit-learn used for?

Ans.The make_circles package in scikit-learn is a utility function that generates a toy dataset of 2D points that form concentric circles with Gaussian noise. This dataset is often used to demonstrate clustering algorithms, including DBSCAN.

Q8.What are local outliers and global outliers, and how do they differ from each other?

Ans.Local outliers are data points that are anomalous with respect to their local neighborhood, while global outliers are anomalous with respect to the entire dataset. Local outliers may not be detected by distance-based methods such as KNN, while global outliers may be easier to detect using such methods.

Q9. How can local outliers be detected using the Local Outlier Factor (LOF) algorithm?

Ans. Local outliers can be detected using the Local Outlier Factor (LOF) algorithm by comparing the density of each data point to the density of its k-nearest neighbors. Data points with a significantly lower density than their neighbors are considered to be local outliers. The key parameter involved in the process is k, which determines the number of nearest neighbors used for density estimation.

Q10. How can global outliers be detected using the Isolation Forest algorithm?

Ans. Global outliers can be detected using the Isolation Forest algorithm by randomly selecting a feature and a split value for each tree in the forest, and then determining the number of splits required to isolate a data point. Data points that require fewer splits to isolate are considered to be more anomalous.

Q11. What are some real-world applications where local outlier detection is more appropriate than global outlier detection, and vice versa?

Ans. Local outlier detection is more appropriate than global outlier detection in situations where the anomalous behavior is localized and not representative of the entire dataset. Examples of such applications include fraud detection in financial transactions and anomaly detection in sensor networks. On the other hand, global outlier detection is more appropriate in situations where the anomalous behavior is spread out across the entire dataset and represents a significant departure from the normal behavior. Examples of such applications include outlier detection in medical studies and network intrusion detection.