

Help on LOF (Local Outlier Factor) method

The local outlier factor is based on a concept of a local density, where locality is given by k nearest neighbors, whose distance is used to estimate the density. By comparing the local density of an object to the local densities of its neighbors, one can identify regions of similar density, and points that have a substantially lower density than their neighbors. These are considered to be outliers. Due to the local approach, LOF is able to identify outliers in a data set that would not be outliers in another area of the data set. For example, a point at a "small" distance to a very dense cluster is an outlier, while a point within a sparse cluster might exhibit similar distances to its neighbors. The local density is estimated by the typical distance at which a point can be "reached" from its neighbors. The definition of "reachability distance" used in LOF is an additional measure to produce more stable results within clusters. Formally, let $k\text{-distance}(A)$ be the distance of the object A to the k -th nearest neighbor. We denote the set of k nearest neighbors as $N_k(A)$. Also define a reachability distance as:

$$\text{reachability} - \text{distance}_k(A, B) = \max[k - \text{distance}(B), d(A, B)]$$

To get more stable results the reachability distance of an object A from B is the true distance of the two objects, but at least the $k - \text{distance}$ of B and Objects that belong to the k nearest neighbors of B (the "core" of B) are considered to be equally distant. The local reachability density of an object A is defined as:

$$\text{lrd}(A) = 1 / \left(\frac{\sum_{B \in N_k(A)} \text{reachability} - \text{distance}_k(A, B)}{|N_k(A)|} \right)$$

$\text{lrd}(A)$ is the inverse of the average reachability distance of the object A from its neighbors i.e., the distance at which A can be "reached" from its neighbors. The local reachability densities are then compared with those of the neighbors using

$$\text{LOF}_k(A) := \frac{\sum_{B \in N_k(A)} \frac{\text{lrd}(B)}{\text{lrd}(A)}}{|N_k(A)|} = \frac{\sum_{B \in N_k(A)} \text{lrd}(B)}{|N_k(A)| \times \text{lrd}(A)}$$

$\text{LOF}_k(A)$ is the average local reachability density of the neighbors divided by the object's own local reachability density. A value of approximately 1 indicates that the object is comparable to its neighbors (and thus not an outlier). A value below 1 indicates a denser region (which would be an inlier), while values significantly larger than 1 indicate outliers.

Screen example with $k = 5$ is depicted below:

