

with zero corresponding to a NA value. Next we want to identify the bivariate outliers by iterating the filter over all possible pairs of variables. Consider a pair of variables $\mathbf{X}^{(jk)} = \{(\mathbf{X}_{ij}, \mathbf{X}_{ik})\}, i = 1, \dots, n$. The initial location and dispersion estimators are, respectively, the coordinate-wise median and the 2×2 sub-matrix $S^{(jk)}$ of the estimate S computed by the generalized S-estimator on non-filtered data. Note that, this ensure the positive definiteness property for S and each $d \times d$ sub-matrix corresponding to a subset of d variables. For bivariate points with no flagged components by the univariate filter we compute the squared Mahalanobis distance $\Delta_i^{(jk)}$ and hence apply the bivariate filter, for all $1 < j < k < p$. At the end we want to identify the cells (i, j) which have to be flagged as cell-wise outliers. The procedure used for this purpose is described in ? and reported here. Let

$$J = \{(i, j, k) : \Delta_i^{(jk)} \text{ is flagged as bivariate outlier}\}$$

be the set of triplets which identifies the pairs of cells flagged by the bivariate filter in rows $i = 1, \dots, n$. For each cell (i, j) in the data, we count the number of flagged pairs in the i -th row in which the considered cell is involved:

$$m_{ij} = \#\{k : (i, j, k) \in J\}.$$

In absence of contamination, m_{ij} follows approximately a binomial distribution $Bin(\sum_{k \neq j} \mathbf{U}_{jk}, \delta)$ where δ represents the overall proportion of cell-wise outliers undetected by the univariate filter. Hence, we flag the cell (i, j) if $m_{ij} > c_{ij}$, where c_{ij} is the 0.99-quantile of $Bin(\sum_{k \neq j} \mathbf{U}_{jk}, \mathbf{0.1})$. Finally, we perform the p -variate filter as described in subsection to the full data matrix. Detected observations (rows) are directly flagged as p -variate (case-wise) outliers. We denote the procedure based on univariate, bivariate and p -variate filters, HS-UBPF.

0.1 A sequencing filtering procedure

Suppose we would like to apply a sequence of k filters with different dimension $1 \leq d_1 \leq d_2 \leq \dots \leq d_k \leq p$. For each d_i , $i = 1, \dots, k$, the filter updates the data matrix adding NA values to the d_i -tuples identified as d_i -variate outliers. In this way, each filter applies only those d_i -tuples that have not been flagged as outliers by the filters with lower dimension. Initial values for each procedures rather than d_1 would be obtained by applying the GSE to the actual filtered values. This procedure aims to be a valid alternative to that used in the presented HS-UBPF filter to perform a sequence of filters with different dimensions. However, this is a preliminary idea, indeed it has not been implemented yet.