# Copula-based anomaly scoring and localization for large-scale, high-dimensional continuous data

GÁBOR HORVÁTH, Dept. of Networked Systems and Services, Budapest University of Technology and Economics, Hungary

EDITH KOVÁCS, Dept. of Differential Equations, Budapest University of Technology and Economics
Faculty of Informatics, University of Debrecen, Hungary

ROLAND MOLONTAY, Dept. of Stochastics, Budapest University of Technology and Economics,
MTA-BME Stochastics Research Group
Faculty of Informatics, University of Debrecen, Hungary

SZABOLCS NOVÁCZKI, Nokia, Bell Labs, Hungary

The anomaly detection method presented by this paper has a special feature: it does not only indicate whether an observation is anomalous or not but also tells what exactly makes an anomalous observation unusual. Hence, it provides support to localize the reason of the anomaly.

The proposed approach is model-based; it relies on the multivariate probability distribution associated with the observations. Since the rare events are present in the tails of the probability distributions, we use copula functions, that are able to model the fat-tailed distributions well. The presented procedure scales well; it can cope with a large number of high-dimensional samples. Furthermore, our procedure can cope with missing values, too, which occur frequently in high-dimensional data sets.

In the second part of the paper, we demonstrate the usability of the method through a case study, where we analyze a large data set consisting of the performance counters of a real mobile telecommunication network. Since such networks are complex systems, the signs of sub-optimal operation can remain hidden for a potentially long time. With the proposed procedure, many such hidden issues can be isolated and indicated to the network operator.

CCS Concepts: • **Computing methodologies** → **Anomaly detection**; *Unsupervised learning*; *Dimensionality reduction and manifold learning*; • **Networks** → *Mobile networks*;

Additional Key Words and Phrases: Anomaly scoring, unsupervised learning, copula fitting

Authors' addresses: Gábor Horváth, Dept. of Networked Systems and Services, Budapest University of Technology and Economics, Magyar Tudósok krt. 2, Budapest, Hungary, 1117, ghorvath@hit.bme.hu; Edith Kovács, Dept. of Differential Equations, Budapest University of Technology and Economics
Faculty of Informatics, University of Debrecen, Műegyetem rkp. 3. Budapest, Hungary, 1111, kovacsea@math.bme.hu; Roland Molontay, Dept. of Stochastics, Budapest University of Technology and Economics,
MTA-BME Stochastics Research Group
Faculty of Informatics, University of Debrecen, P.O. Box 91, Budapest, Hungary, 1521, molontay@math.bme.hu; Szabolcs Nováczki, Nokia, Bell Labs, Bókay János utca 36-42, Budapest, Hungary, 1083, szabolcs.novaczki@nokia-bell-labs.com.

# 1 INTRODUCTION

Anomaly detection refers to the process of identifying unexpected objects or patterns, which do not conform to the usual behaviour. On the other hand, there is no canonical definition of anomaly, generally an object is called anomaly if it is different from normal instances with respect to its features and it is rare in the dataset [26]. In the literature, these instances are also often referred to as outliers, rare/extreme events, discordant objects [62]. The detection of "not-normal" observations has attracted a lot of research interest from the machine learning community since it has a wide variety of practical applications, including network intrusion, credit card fraud, health anomaly detection and so on [12].

We can distinguish between three main types of anomaly detection setups based on the availability of labels in the dataset [26]. Supervised anomaly detection means that the data is fully labelled and an ordinary classifier can be used after dealing with the unbalanced class distribution [18]. For semi-supervised anomaly detection, the training data consists of normal instances without any anomalies, one-class classification algorithms can be used [39], moreover, density estimation methods can be also used to model the density function of the normal class [36, 37], then the instances that do not conform to the normal profile are identified as anomalies. Unsupervised anomaly detection is performed on unlabeled data, taking only intrinsic properties of the dataset into account [26].

Anomaly detection algorithms can be further classified based on the specific data types and domains that they are suitable for, such as time series, categorical attributes, item sets, graphs, spatial data [26, 65]. The present study focuses on unsupervised anomaly detection for numerical high-dimensional data in Euclidean space. Due to the "curse of dimensionality" and the presence of irrelevant attributes high-dimensional anomaly detection is particularly challenging [65]. In this paper, we propose a probabilistic model-based unsupervised anomaly detection method that works well with a large number of high-dimensional numerical observations. By numerical data we mean that all variables are non-categorical and have a continuous distribution. This paper also describes a possible way to model the joint probability density based on a general real-valued multidimensional dataset. The proposed modelling approach can be used to assign anomaly scores to any point of the multidimensional space: observations from low density areas of the space receive high anomaly scores and vice versa.

Throughout the years, several algorithms have been proposed for unsupervised anomaly detection in high-dimensional numerical data, for extensive reviews and comparative studies we refer the reader to [3, 26, 62, 65].

A high number of outlier detection approaches rely on the concept of neighborhood, either using $k$-nearest neighbor-distances [4] or using a density based approach, i.e. comparing the number of instances in $\varepsilon$-neighborhood of the object to the $\varepsilon$-neighborhood of the object's neighbors, e.g. Local Outlier Factor (LOF) [10]. Although various variants of LOF have been proposed, it is important to note that neighborhood is less meaningful if the dimension is high [65]. A possible approach to handle high-dimensionality is to use dimension reduction methods to improve outlier detection [48]. The assumption of dimension reduction based outlier detection methods is that a single subspace is sufficient to identify anomalies.

Another related approach is using ensemble learning methods and combine different subspaces for anomaly detection. The idea of "feature bagging" method is to derive outliers in randomly selected feature subsets using multiple outlier detection algorithms and combine the outputs to achieve more stable and effective results by avoiding the curse of high dimensionality [38, 47]. A sequential ensemble-based framework was proposed to mutually refine feature selection and outlier scoring [51].

Subspace outlier detection techniques define outlierness with respect to specific subspaces based on the observation that outliers are often embedded in locally relevant subspaces but in the full-dimensional space they are covered by irrelevant attributes [2, 24]. The challenge is to identify the relevant subspaces. Several techniques have been proposed for subspace identification recently, including rarity based techniques [63], density-based and grid-based subspace clustering approaches [33, 45]. Subspace outlier detection (SOD) is a method that works without a previous clustering step [32]. Another approach is to search for correlated subspaces of features [49, 50], it was also proved that the subspace search problem can be transformed into a problem of clique mining for highly correlated features. In a recent paper a new subspace analysis approach was proposed named Agglomerative Attribute Grouping (AAG) that relies on a novel multi-attribute information theoretical measure which evaluates the "information distance" between groups of features [5].

Detecting outlying subspaces is also advantageous for the purpose of outlier interpretation, i.e. identifying the attributes that contribute the most to abnormality [61]. Other methods have been also proposed to address the issue of outlier interpretation, such as Local Outlier Detection with Interpretation (LODI) [17] and Contextual Outlier INterpretation (COIN) [43].

Based on the observation that angles between data points are more meaningful in high dimension than distances, angle-based outlier detection methods have also been proposed [34]. Angle-based methods rely on the assumption that the variance of angles between an outlier and other data points is lower than for normal observations, since for normal data points, other objects are distributed in all direction, while for outliers, they are more concentrated in some direction. A random projection based efficient approximation variant was proposed in [52].

Probabilistic anomaly detection methods have also been extensively used for high-dimensional data. The general scheme is to estimate the density function of a dataset $X$ by calibrating the model parameters $\theta$ and identifying outliers as the observations having the smallest likelihood $\mathbb{P}(X|\theta)$. The density function can be estimated by Gaussian Mixture Model using Expectation-Maximization algorithm [59], by Dirichlet Process Mixture Model [21], by Kernel Densisty Estimators [59], and by Robust Kernel Density Estimator [29] among others.

Isolation methods are also quite widespread, isolation forest compute an anomaly score using random forest [40–42]. Instances which are easy to isolate, i.e., have short average path lengths on the isolation trees, are considered anomalous. To reduce the dimension of the feature space, only a subset of the features is considered based on their kurtosis, since it is observed that kurtosis is sensitive to the presence of anomalies.

Neural networks are also used for anomaly detection in high-dimensional data. A reconstruction-based self-organizing neural network, called Grow When Required [44] can be applied for outlier detection [18], furthermore, Recurrent Neural Networks have also been used for network traffic anomaly detection [54]. A Pattern Anomaly Value (PAV) [13] based anomaly detection method for sensor data was introduced in [20].

We can observe that a high number of unsupervised anomaly detection methods have been proposed throughout the years. On the other hand, it is important to note that there is no canonical method, mainly because it is really difficult to interpret and compare various methods. The continuous outlier scores provided by most of the algorithms vary widely in in their scale, range, and meaning [65]. Evaluation of unsupervised anomaly detection methods is a notorious problem: efficiency, effectiveness, interpretability, scalability, memory consumption and robustness should also be considered. For some possible evaluation measures, strategies and empirical evaluation of various anomaly detection techniques on benchmarks datasets we refer to [11, 62, 65].

In this paper, we introduce a method, which tackles the high dimensional feature space by finding some special two dimensional subspaces, which also aims to minimize redundancy. An

other important benefit of the presented approach is that it does not only give an anomaly score but also indicates why exactly an observation is anomalous, making it easier to find and fix the source of the problem. Furthermore, the proposed method scales well for high dimension and for large sample size, and is able to return anomaly scores for samples with missing data as well.

## 2 OUTLINE OF THE PROCEDURE

Our procedure is based on probabilistic modeling. The main idea is to determine the joint distribution of the random variables corresponding to the features, and to assign anomaly scores to the observations based on the density of the observations. The subspaces where an observation deviates from the majority of the samples indicate which variables are involved in the anomaly. However, the realization of this high-level description of the algorithm faces some challenges that the proposed approach aims to overcome:

- Obtaining the joint distribution of a high-dimensional dataset is difficult. Simple solutions for this problem aim to fit the dataset by a multivariate Gauss distribution or by a mixture of multivariate Gaussians. However, Gaussian density functions have a fast decay, they fail to capture the heavy-tailed behavior where most anomalies take place, making them inappropriate to use for anomaly detection.
- Even if the joint density function of the dataset is available, translating the density (a number between zero and infinity) to an anomaly score (between zero and one) having a physical interpretation, is not straightforward.
- Providing aid to localize the reason of the anomaly in a high-dimensional space, such that a human can interpret it and take the necessary actions, is also challenging.

In our approach (see Figure 1), the high-dimensional problem is decomposed into several two-dimensional ones. As the first step, we identify the most relevant *pairs of variables* whose joint behavior retains as much information from the high-dimensional joint behavior as possible. The details of this step are provided in Section 3.

In the next step, described by Section 4, the joint distributions of the selected pairs of variables are characterized by a copula-based method.

Finally, the anomaly scores are calculated from the joint densities, as presented by Section 5. For each observation, our procedure is able to report both an overall anomaly score, and individual anomaly scores for the selected variable pairs. Knowing the variable pairs affected helps to localize the reason of the anomaly.

## 3 DECOMPOSITION OF THE HIGH DIMENSIONAL SPACE

The direct modeling of the joint behavior of high-dimensional observations with reasonable accuracy is intractable. To cope with the exploding complexity, we use only some of the two-dimensional marginal probability distributions in the presented method, those that are the most relevant from the anomaly detection point of view.

The question that naturally arises is how to chose a relatively small number of pairs from all possible pairs of random variables, while retaining as much information as possible. The answer is inspired by the work of Chow and Liu [15], where a similar problem has been solved for discrete joint probability distributions. In this section, we provide a summary on this method, and adapt it to the case of continuous variables.

Let us consider $\mathbf{X} = (X_1, \ldots, X_d)^T$ a $d$-dimensional discrete random vector with probability distribution $P(\mathbf{X})$. The probability of a realization $\mathbf{x}$ is denoted by $P(\mathbf{x})$, and the set of indices of the random variables are denoted by $V = \{1, \ldots, d\}$. Furthermore, we denote the probability distributions of the two-dimensional marginals by $P(X_i, X_j), i, j \in V, i \neq j$.
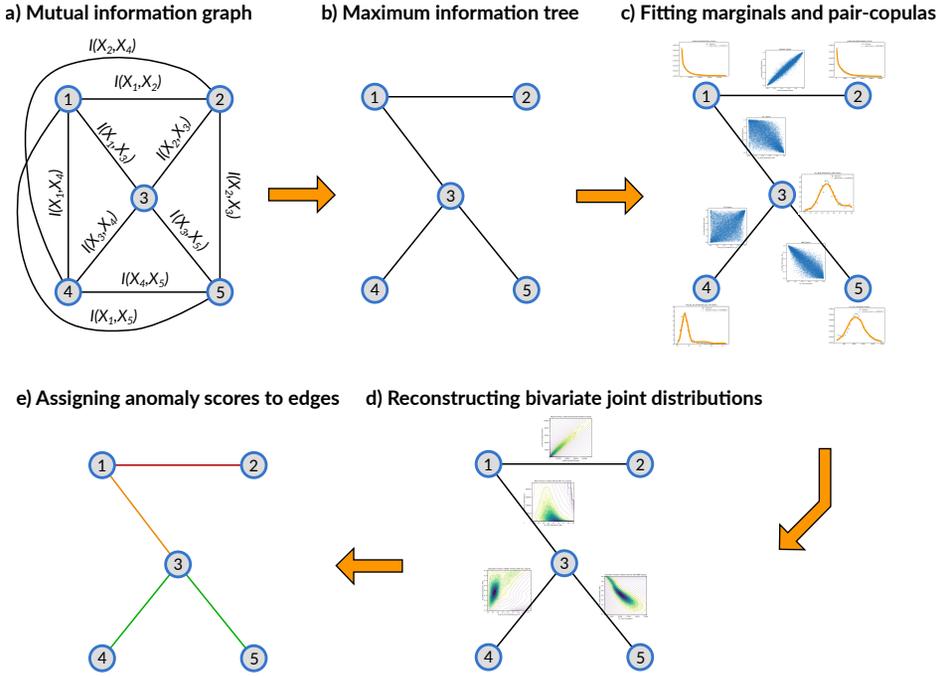
Fig. 1. Concept of the presented anomaly scoring method

In [15], a product-form probability distribution has been introduced, associated with a spanning tree over the set of indices, defined by

$$P_{Ch\_L}(\mathbf{X}) = \frac{\prod\limits_{(i,j)\in T} P(X_i, X_j)}{\prod\limits_{i \in V} [P(X_i)]^{v_i-1}},$$ (1)

where $T$ is a spanning tree defined on the vertex set $V = \{1, \ldots, d\}$, and $v_i$ denotes the number of such two-dimensional marginal probability distributions in the numerator of (1) which contain variable $X_i$. This type of probability distribution has the property that its two-dimensional marginals present in its formula coincide with the two-dimensional marginals of $P(\mathbf{X})$.

Our aim is to find a probability distribution $P_{Ch\_L}(\mathbf{X})$ of (1) to approximate the probability distribution $P(\mathbf{X}) = (X_1, \ldots, X_d)^T$, with the goodness of fitting quantified by the Kullback-Leibler divergence. Hence, the problem can be formulated as searching for a probability distribution $P_{Ch\_L}(\mathbf{X})$ which minimizes the Kullback- Leibler divergence between $P(\mathbf{X})$ and $P_{Ch\_L}(\mathbf{X})$.

To solve this problem, [15] has introduced an undirected complete graph defined on the set of vertices $V$. The edges of the graph are weighted by the mutual information of the two dimensional marginal probability distributions corresponding to the two vertices connected. In that paper, it has been proven that the Kullback-Leibler divergence between $P(\mathbf{X})$ and $P_{Ch\_L}(\mathbf{X})$ can be expressed as

$$KL(P_{Ch\_L}(\mathbf{X}), P(\mathbf{X})) = \sum_{i=1}^{d} H(X_i) - \sum_{(i,j)\in T} I(X_i, X_j) - H(\mathbf{X}),$$

where $I(X_i, X_j)$ is the mutual information defined by

$$I(X_i, X_j) = H(X_i) + H(X_j) - H(X_i, X_j), \tag{2}$$

and $H(\cdot)$ denotes the entropy of a random variable or a random vector.

It is easy to see that $KL(P_{Ch\_L}(\mathbf{X}), P(\mathbf{X}))$ is minimal when we take edges from the complete graph along the spanning tree having the maximum weight [15], i.e.,

$$T^* = \underset{T}{\operatorname{argmax}} \left( \sum_{(i,j) \in T} I(X_i, X_j) \right).$$

The maximum information tree $T^*$ can be obtained by applying Prim's or Kruskall's algorithm.

It is important to mention that the spanning tree structure encodes conditional independencies between the random variables associated to the vertices. If two nodes are not connected, the corresponding random variables are conditionally independent given any random variable on the path between them. Therefore, the best fitting $P_{Ch\_L}(\mathbf{X})$ exploits the conditional independencies existing in the multivariate probability distribution that reduces the redundancy between the corresponding random variables.

In many practical cases, such as in the case described in Section 6, the samples are obtained from a continuous probability distribution instead of a discrete one. In such cases the samples can be transformed to discrete, and the idea of Chow and Liu can still be applied. For this transformation we use a partitioning such that each interval contains the same number of realizations for all random variables. This way we get $H(X_i) = H(X_j)$, for all $i, j \in \{1, \ldots, d\}, i \neq j$, and therefore we have that $I(X_i, X_j)$ depends only on $H(X_i, X_j)$ (see (2)).

Since the maximum of $\sum_{i,j \in T} I(X_i, X_j)$ is obtained when $\sum_{i,j \in T} H(X_i, X_j)$ is minimal, we weight the edges of the complete graph by the entropy of the bivariate random variables connected, and construct the minimal weighted spanning tree of this graph. As a result, we obtain a set of $d - 1$ variable pairs, whose bivariate marginal probability distributions capture as much information on the joint behavior as possible (Step b. in Figure 1).

There are software packages available to obtain the maximum information tree, also for the case of continuous variables [55]. These packages use Kendall's tau instead of the mutual information for weighting the graph. However, for big data the computation of Kendall's tau is computationally expensive, therefore we recommend to use the mutual information on the discretized data for this purpose. The tree structure is not sensitive for the chosen metric, the obtained three structure provided by these methods was the same in all of the studied cases.

## 4  COPULA-BASED MODELING OF PAIRS OF VARIABLES

The proposed method relies on *copulas* to model the joint distributions of the variable pairs obtained by the maximum information tree. Copula-based methods are widely used in modelling joint probability distributions (mainly in the field of financial mathematics) due to their benefits in representing heavy tails [14]. Since anomalies can be considered as events falling in the tail of the probability distribution, copulas are very useful in modeling anomalies.

According to Sklar's theorem, the joint distributions are uniquely decomposed to the marginal distributions and a so called copula function, that characterizes only the dependency structure between the variables, and is independent from the marginal distributions. Consequently, copulas enable the separate modeling of the dependence structure and the marginals; we are going to rely heavily on this feature in the presented method.

## 4.1 A short overview on copulas

Before detailing the exact kind of copulas used in this paper, let us review some concepts related to copulas (for more details see [46]).

*Definition 4.1.* A function $C : [0; 1]^d \rightarrow [0; 1]$ is called a $d$-dimensional copula if it satisfies the following conditions:

(1) $C(u_1, \ldots, u_d)$ is monotonously increasing in each component $u_i$,
(2) $C(u_1, \ldots, u_{i-1}, 0, u_{i+1}, \ldots, u_d) = 0$ for all $u_k \in [0; 1]$, $k \neq i$, $i = 1, \ldots, d$,
(3) $C(1, \ldots, 1, u_i, 1, \ldots, 1) = u_i$ for all $u_i \in [0; 1]$, $i = 1, \ldots, d$,
(4) $C$ is $d$-increasing, i.e for all $(u_{1,1}, \ldots, u_{1,d})$ and $(u_{2,1}, \ldots, u_{2,d})$ in $[0; 1]^d$ with $u_{1,i} < u_{2,i}$ for all i, we have

$$\sum_{i_1=1}^{2} \cdots \sum_{i_d=1}^{2} (-1)^{\sum_{j=1}^{d} i_j} C\left(u_{i_1,1}, \ldots, u_{i_d,d}\right) \geq 0.$$

Thus, copulas are $d$ dimensional distribution functions such that each of their marginals are uniformly distributed.

Due to Sklar's theorem [57], if $X_1, \ldots, X_d$ are continuous random variables defined on a common probability space, with univariate marginal cdf's $F_{X_i}(x_i)$ and joint cdf $F_{X_1, \ldots, X_d}(x_1, \ldots, x_d)$, there exists a unique copula function $C_{X_1, \ldots, X_d}(u_1, \ldots, u_d) : [0; 1]^d \rightarrow [0; 1]$, such that by the substitution $u_i = F_i(x_i)$, $i = 1, \ldots, d$, we get

$$F_{X_1, \ldots, X_d}(x_1, \ldots, x_d) = C_{X_1, \ldots, X_d}(F_1(x_1), \ldots, F_d(x_d)), \text{ for all } (x_1, \ldots, x_d)^T \in R^d. \qquad (3)$$

This relation is of principal importance. It states that *every* multivariate joint distribution can be fully reconstructed based on its marginals and its copula.

The density function $f_{x_1, \ldots, x_d}(x_1, \ldots, x_d)$ can be expressed in the following way:

$$f_{x_1, \ldots, x_d}(x_1, \ldots, x_d) = \frac{\partial^d F(x_1, \ldots, x_d)}{\partial x_1 \ldots \partial x_d} = \left. \frac{\partial^d C(u_1, \ldots, u_d)}{\partial u_1 \ldots \partial u_d} \right|_{u_i = F(x_i)} \prod_{i=1}^{d} \frac{\partial F(x_i)}{\partial x_i}$$

$$= c(F_1(x_1), \ldots, F_d(x_d)) \prod_{i=1}^{d} f(x_i), \qquad (4)$$

where $c(u_1, \ldots, u_d)$ is the copula density function.

Modeling high-dimensional probability distributions by using only a single high-dimensional copula to characterize all the underlying dependencies is possible, although rather restrictive, as there is a high chance that none of the currently known copula families is flexible enough to approximate the complex dependency structure of a real dataset with reasonable accuracy (as shown in [1]). A new approach has been introduced by Joe [28], the so called *pair-copula construction*, that makes possible to capture different kinds of dependencies between the pairs of the variables in the same multivariate probability distribution. In this approach, a copula is expressed as a product of different types of bivariate copulas and bivariate conditional copulas. A useful modeling tool, following the same direction, is the so called R-vine structure, that was introduced by Bedford and Cooke [6, 7] and described in more details by Kurowicka and Cooke [35]. The drawback of using R-vine copulas in higher dimensions is that the number of parameters grows very fast with the dimensions. To address this issue, a special kind of vine-copula was proposed, called truncated vine-copula. This concept has been introduced by Brechmann [8] and a more general version has been presented in [30, 31].

We do not review the R-vine graph structures here, as we need only a special case of them, the truncated vines [8, 9]. In our approach, we will use the so called truncated vine copula truncated at level one, that is directly related to Chow Liu probability distributions (see Section 3). Let us denote the bivariate probability distribution function of $X_i, X_j$ by $F_{i,j}$, the bivariate probability density function by $f_{i,j}$, and the bivariate copula density function corresponding to $f_{i,j}$ by $c_{i,j}$, for $i, j \in V$.

Let us consider the approximation of the joint probability density function $f_{X_1,\ldots X_d}(\mathbf{x})$ given by the by a Chow-Liu type density function:

$$f_{Ch\_L}(\mathbf{x}) = \prod_{e \in T, s \in V} \frac{f_{j(e)k(e)}\left(x_{j(e)}, x_{k(e)}\right)}{f_s(x_s)^{v_s - 1}}, \tag{5}$$

where $v_s$=the number of edges connected to vertex $s$ (see (1) for the discrete case).

The probability density function $f_{Ch\_L}(\mathbf{x})$ can be expressed as follows by dividing and multiplying with the marginal probability densities of the bivariate densities in the numerator, and then by applying (4) to the bivariate case:

$$\begin{aligned} f_{Ch\_L}(\mathbf{x}) &= \prod_{i=1}^{d} f(x_i) \prod_{e \in T} \frac{f_{j(e)k(e)}\left(x_{j(e)}, x_{k(e)}\right)}{f_{j(e)}\left(x_{j(e)}\right) \cdot f_{k(e)}\left(x_{k(e)}\right)} = \\ &= \prod_{i=1}^{d} f(x_i) \prod_{e \in T} c_{j(e)k(e)}\left(F_{j(e)}\left(x_{j(e)}\right), F_{k(e)}\left(x_{k(e)}\right)\right). \end{aligned} \tag{6}$$

It is easy to see that from (4) and from (6) it follows that $\prod_{e \in T} c_{j(e)k(e)}\left(F_{j(e)}\left(x_{j(e)}\right), F_{k(e)}\left(x_{k(e)}\right)\right)$ is a valid copula density.

Observe that in (6) only pair-copulas are involved, this is the reason why these constructions are also called pair-copula constructions. Formula (6) also shows that the univariate marginal probability distributions and the bivariate copulas can be fitted separately.

Using the elements introduced above, the inference of the multivariate density function is boiled down to the following steps (see also Figure 1):

- The selection of a specific tree-structure. For this purpose we use the maximum information tree described in Section 3.
- The choice of pair-copulas involved in formula (6). For each edge in the maximum information tree, thus for each variable pair, a pair-copula is constructed and its parameters are estimated.
- The modeling of the univariate marginals. For every vertex of the maximum information tree, thus for every variable, the marginal distribution is estimated with an appropriately chosen distribution family.
- Construction of the two-dimensional joint distributions. For every edge of the maximum information tree, the joint distribution of the variable pair is reconstructed from the two related marginals and the copula.

Let us consider i.i.d observations $D = \left\{\mathbf{x}^1, \ldots, \mathbf{x}^N\right\}, \mathbf{x}^m \in \mathbb{R}^d$. Now we fit a Chow-Liu density function to the data. The expression of the loglikelihood function of Chow-Liu density function (6) fitted to the data $D$ is the following:

$$\ln L\left(f_{Ch_L}\left(D\right),\alpha_{lm},\beta\right) = \sum_{e\in T}\sum_{m=1}^{N}\ln c_{j(e)k(e)}\left(F_{j(e)}\left(x_{j(e)}^m\right), F_{k(e)}\left(x_{k(e)}^m\right),\alpha_e\right) +$$

$$+ \sum_{m=1}^{N}\ln f\left(x_i^m,\beta\right), \tag{7}$$

where $\alpha_e, \beta$ are the parameters of the pair copulas in the edges, and the marginal density functions respectively. From (7) it follows that the loglikelihood of the Chow-Liu approximation can be expressed as the sum of loglikelihoods of the pair-copulas and the sum of the loglikelihood of the marginal pdf's. This makes our fitting tractable, paralelizable and flexible to multiple kind of dependencies.

The following three subsections provide solutions for the last three items, for the fitting of pair-copulas, the univariate marginals, and the construction of the joint distribution.

## 4.2 Fitting the pair-copulas

For every variable pair in the maximum information tree, our procedure first extracts the empirical copula. To do so, both variables need to be transformed to pseudo-observations, hence they are replaced by the normalized ranked data. The pseudo-observations corresponding to $N$ realizations of a $d$-dimensional random vector $x_i = (x_{i1},\ldots,x_{id})^T$, $i\in\{1,\ldots,N\}$ are defined by $u_{ij} = \frac{r_{ij}}{N+1}$, $i\in\{1,\ldots,N\}, j\in\{1,\ldots,d\},$, where $r_{ij}$ denotes the rank of $x_{ij}$ among all $x_{kj}, k\in\{1,\ldots,N\}$ [55]. This way we get a two-dimensional data series where both variables are uniformly distributed individually, and their joint values are the realizations of their copula distribution.

To find the appropriate copula model for the empirical copula, we consider a set of copula families, perform fitting the observations with all of them, and select the one providing the highest likelihood.

In our method the following copula families (and their rotations) are considered:

- Archimedean copulas, including the Gumbel, Clayton, Joe, Frank, BB1, BB6, BB7, BB8 copulas. These copulas have a low number of parameters and an explicit copula density, that makes the fitting process relatively efficient.
- Elliptical copulas [22], including the Gaussian and the Student-t copula. Contrary to the Archimedean copulas the elliptical ones are symmetric. They do not have an explicit density, hence their fitting does not scale very well with the amount of training data.
- Non-parametric copula: the TTL copula. It makes use of a local-likelihood transformation kernel density estimation (see [25]). This method has a large number of parameters (900) and yields a less smooth density; in the presented procedure it serves as a remedy when the empirical copula is so asymmetric that the above listed copula families fail to fit it.

For this step of the procedure we were relying on the built-in copula fitting capabilities of the vinecopulib package[1].

## 4.3 Fitting the univariate marginal distributions

The concept of fitting the marginals for each variable (a vertex in the maximum information tree) is the same as in case of copulas: having defined a wide set of distribution families we repeat the fitting with all of them and choose the one that has the highest likelihood for the training data.

The following univariate distribution families were the candidates for fitting the marginals:

---

[1]https://github.com/vinecopulib/vinecopulib

- Simple univariate distributions like the exponential, student-t, log-normal, and the Pareto distributions. Finding the maximum likelihood estimation for these distributions is simple, but they often do not have the flexibility to fit real-word dataset accurately enough.
- Gaussian mixture models. Gaussian mixture models are widely used for fitting real datasets. With the expectation-maximization algorithm it is possible to find its optimal parameters rather efficiently. However, the density function of Gaussian mixtures has a fast decay, they are not the optimal choice for distributions having a heavier tail.
- Hyper-Erlang distributions. Hyper-Erlang distributions are the mixtures of Erlang distributions. Efficient expectation-maximization algorithm for the inference is available for these distributions, too (see [60]). While the gamma distributions are more general than the Erlang distributions, interestingly we got far better fitting results with hyper-Erlangs. In fact, they turned out to be the most flexible choices for the majority of the variables in our use case in Section 6.
- Mixture of beta distributions. To model distributions with finite support, the mixture of beta distributions provide an appealing option. We have successfully applied the hybrid EM – moment fitting method published in [56].

## 4.4 Construction of the joint density functions

Once the univariate marginals and the pair-copulas are available, the two-dimensional joint densitiy functions can be reconstructed based on (4), yielding

$$f_{X_i, X_j}(x_i, x_j) = c_{X_i, X_j}(F_{X_i}(x_i), F_{X_j}(x_j)) \cdot f_{X_i}(x_i) \cdot f_{X_j}(x_j), \tag{8}$$

and the overall density of a $d$ dimensional observation is given in product form by (6).

## 5 ANOMALY SCORES DERIVED FROM THE JOINT PROBABILITY DISTRIBUTION

Let us consider a $d$-dimensional random vector $X = (X_1, \ldots X_d)^T$ taking values in $\chi \subset \mathbb{R}^d$, with a bounded probability density function denoted by $f(\mathbf{x})$. The aim of this section is to define an anomaly scoring function which is defined on the $d$-dimensional support of the random vector $\mathbf{X}$ and takes values in the $[0; 1]$ interval.

There are several distributions for which it is easy to obtain an anomaly score from the density functions. In the case of a symmetrical probability distribution, such as the univariate Gaussian distribution, it is known that the realizations which deviate from the mean value more than $2\sigma$ are considered rare. The probability that a realization of a Gaussian random variable takes place in the set $A_{0.9545}^G = \{x \mid |x - \mu| < 2\sigma\}$ is 0.9545. For non-symmetrical, uni-modal distributions, a plausible set $A_{0.9545}$ would contain values between the 2.275 and 97.725 percentiles.

For multivariate distributions with multimodal density functions, however, it is not obvious how to define such an anomaly scoring function. Results related to this topic were developed in [19, 53], based on the idea of splitting the feature space in two halves, and observations not falling into *the minimum volume set* (see Definition 5 in [23]) are considered to be anomalous. In our approach we do not only classify the realizations as normal or anomalous but we also assign anomaly scores to them. Other recent work in this direction can be found in [16], where a so called *Mass Volume Curve* (MV curve) is defined. The idea behind MV curves is to define an ordering on $\chi \subset \mathbb{R}^d$ according to which realizations of $\mathbf{X}$ can be ranked related to how rare they are.

In this paper we rely on the *Mass Volume Curve*, too. In this section we review the main elements of this approach and provide and algorithm to obtain the anomaly scores of any high-dimensional observations based on the joint density function efficiently.

Let us start with the definition of the scoring function and the level set.
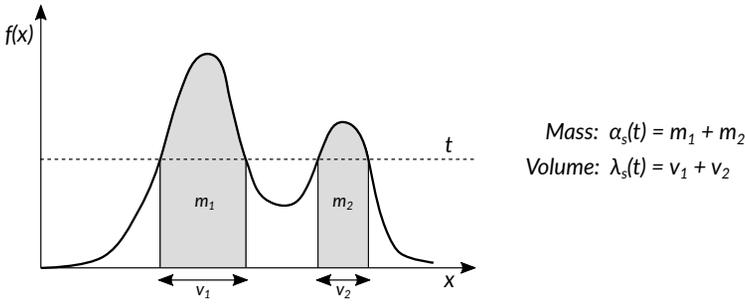
Fig. 2. Mass and volume on a bi-modal density function

*Definition 5.1 (from [16]).* A *scoring function* is any measurable function $s : \chi \to \mathbb{R}_+$ that is integrable with respect to the Lesbesque measure.

A *level set* corresponding to a scoring function $s$ and a level $t$ is given by $A_t = \{\mathbf{x} \in \chi \mid s(\mathbf{x}) \geq t\}, t \in [0, \infty]$ [2].

Note that the "scoring function" as defined above is not the same as our anomaly scoring function defined later. It is easy to see that if $t \geq t'$ then $A_t \subset A_{t'}$. Next, we review the two most important notions, the *mass* and the *volume* of a level set.

*Definition 5.2.* For a scoring function $s$ and a level $t$, $\alpha_s(t) = P(s(\mathbf{X}) \geq t)$ is called the mass of the level set $A_t$.

*Definition 5.3.* For a scoring function $s$ and a level $t$, $\lambda_s(t) = \lambda(\{\mathbf{x} \in \chi \mid s(\mathbf{x}) \geq t\})$ is called the volume corresponding to a level $t$, with respect to the Lesbesque measure.

Based on [53], we consider the density function of $\mathbf{X}$ as a scoring function. In this case $\lambda_f(t)$ is the contour cluster of the density function $f(x)$ at level $t$.

Mass, volume, and their relation with anomaly scoring can be interpreted easier though an example. In Figure 2 there is a bi-modal univariate density function, and a level $t$. The mass corresponding to level $t$ is the integral of the density function where it is greater than level $t$, and the volume covers the intervals on the $x$ axis where the density function is above level $t$. Hence, the volume covers the most dense part of the density function and the mass is the probability that a realization falls into this part.

The problem of anomaly detection can be formulated via the so called minimum volume sets. Minimum volume set is a concept introduced in [19, 53] as a solution of the following constrained minimization problem

$$A_\alpha^* = \arg\min_{A \subset \chi} \lambda(A), \text{ such that } P(\mathbf{X} \in A) \geq \alpha, \tag{9}$$

thus, we are looking for the smallest possible volume yielding mass $\alpha$. In [16] it is stated that if the probability density function is bounded and has no flat parts, for any $\alpha \in ]0; 1[$ there exists a unique minimum volume set $A_\alpha^*$ whose mass is equal to $\alpha$. The "rare realizations" are those that belong to the complementary set, i.e., $\chi \backslash A_\alpha^*$, when $\alpha$ is sufficiently large (close to 1).

After these preliminaries we can define the theoretical anomaly score as follows.

*Definition 5.4.* The theoretical anomaly score of a realization $\mathbf{x}$ is defined by

$$\mathcal{A}(\mathbf{x}) = \inf\{\alpha : \mathbf{x} \in A_\alpha^*\}. \tag{10}$$
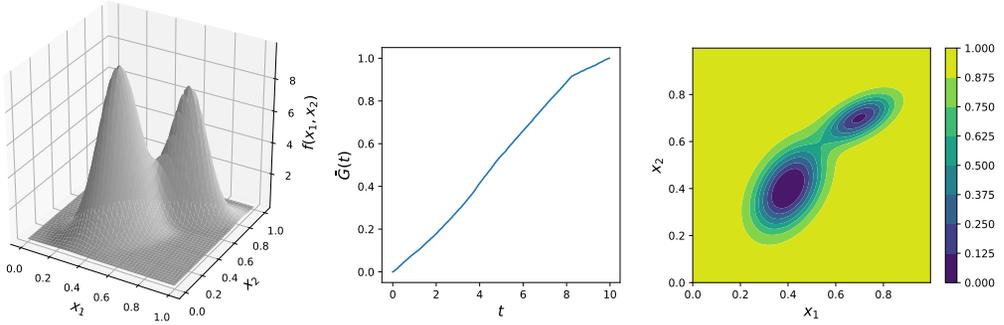
—————
[2]In [16] $t \in [-\infty, \infty]$

Fig. 3. Obtaining anomaly score based on the pdf, in case of a 2-dimensional Gaussian mixture distribution. Left-most: the density, in the middle: $\bar{G}(t)$, right-most: the anomaly scores of the two-dimensional realizations

It is easy to see that $\mathcal{A}(\mathbf{x}) \in [0; 1]$, and the larger the value of $\mathcal{A}(\mathbf{x})$ is the more rare the realization $\mathbf{x}$ is. The main problem is, however, the numerical computation of $A_\alpha^*$ involving the solution of (9).

To overcome this problem we will now define an approximate anomaly score. To do so we generate i.i.d. random samples from $f(\mathbf{x})$, denoted by $\tilde{\mathbf{x}}^i$, $i = 1, \ldots, m$, where $m$ is the sample size. To each vector $\tilde{\mathbf{x}}^i$ we assign its density $f(\tilde{\mathbf{x}}^i)$. This way we get a univariate discrete random variable $\tau$, which takes values $f(\tilde{\mathbf{x}}^i)$ with probability $p_i = \frac{1}{m}$. Then, the anomaly score is approximated by the complementary distribution function of this random variable (denoted by $\bar{G}(t) = P(\tau > t)$) at $t = f(\tilde{\mathbf{x}}^i)$.

*Definition 5.5.* The approximated anomaly score of any realization $\mathbf{x}$ is defined by

$$\hat{\mathcal{A}}(\mathbf{x}) = \bar{G}(f(\mathbf{x})) = \frac{1}{m} \sum 1_{\{\tau > f(\mathbf{x})\}}. \tag{11}$$

It is clear that the approximated anomaly score takes values in $[0; 1]$, and has the property to be large for the rare events. An example for the anomaly score corresponding to a two-dimensional Gaussian mixture is illustrated in Figure 3.

This procedure can be used to identify anomalous cases based on some of the bivariate projections of the multivariate pdf, and also enables the localization of the anomaly, by highlighting those variables whose relation is unusual.

Whenever we have an incomplete realizations, i.e. the observed vector contains missing values, some links of the tree are missing, but all the others can be interpreted as in the case without the missing values. Hence, we can still obtain anomaly scores based on the unusual relation of the existing values.

## 6 APPLICATION OF THE NEW ANOMALY SCORING METHOD TO REAL MOBILE NETWORK DATA

In this section we illustrate the applicability of the proposed procedure on a real dataset originating from the live commercial LTE network of a Western European mobile operator. Mobile networks are complex systems with many parameters, where discovering that the system performance is sub-optimal is difficult for a human operator.

For decades, telecommunication operators have relied on network domain experts to report problems that affect the network performance and customer experience, with the aim of performance monitoring tools. The conventional approach of the tools is to pre-select a set of Key Performance Indicators (KPIs), based on human knowledge and experience of the domain experts. The trending

of these KPIs are closely monitored, based on pre-defined single-value or multi-value thresholds, and/or pre-built single variate or multivariate time series profiles. If the KPI values exceed the thresholds/profiles, then alarms will be raised to trigger investigation, mitigation and on-site support. The experts also try to derive Root Causes Analysis (RCA) on these problems to fix as early as they can. All the processes are manual or semi-auto with limited support of tools. The traditional tools are mainly based on human knowledge, and that the complex and hidden rules in the telecommunication network are not easy to fully detect, capture and utilize.

Thus there is a strong need to revolutionize network management with AI/machine learning technology. This is essential as today's communication networks are extremely complex systems consisting of hundreds of thousands network elements organized in cooperating, coexisting and overlapping technology layers. The network elements generate huge amounts of versatile data for performance monitoring, optimization and troubleshooting purposes. It is already quite cumbersome to tackle these tasks with traditional approaches and the effects will be even more emphasized in case of 5G. Automated solutions are needed that are capable to analyze the raw data and draw conclusions, generate actionable insights using AI technology. One important area in this field is the predictive detection of anomaly patterns that appear in the data.

Our dataset comprises a large number of performance counters collected at thousands of base stations, recorded in every hour. The original set of counters has been filtered based on data quality requirements, which resulted in 54 counters and $1,500,000$ observations. These counters represent a wide variety of measurements: related to the PDCP (Packet Data Convergence Protocol), RLC (Radio Link Control), MAC (Media Access Control) and PHY (Physical) layers, to the handovers, the activity of the schedulers, the UL (upload) and DL (download) throughput , delay, UL interference etc. As mentioned earlier in the paper, anomaly detection for such high dimensional data is a challenging problem. In the following subsections we present how our new method performs with this mobile network dataset.

## 6.1 Selecting the relevant KPI pairs

According to Section 3, the anomaly analysis of the high-dimensional problem is decomposed to many computationally tractable two-dimensional problems which can be treated separately. The first step of the procedure is the selection of the appropriate two-dimensional KPI pairs. For this purpose we create a graph where the vertices are the KPIs and the edges are weighted by the mutual information between them. Those KPI pairs are selected, which are lying on the maximum weight spanning tree of the graph. For the particular dataset considered in this example the maximum weighted spanning tree (referred to as the maximum information tree) is depicted in Figure 4.

The two-dimensional behavior of the selected KPI pairs is then modeled by two-dimensional distributions, using bivariate copula families. The copula based approach enables the separation of the fitting problem of the marginal distributions and the fitting problem of the dependence structure between them.

## 6.2 Fitting distributions to univariate marginals

As described in Section 4.3, we defined various distribution families, and fitted the data of the various KPIs by maximizing the likelihood. Among the candidate distributions the one having the highest likelihood has been selected. From the total number of 54 KPIs the log-normal distribution gave the best likelihood in 1 case, the student distribution in 2 cases, the beta mixture in 6 cases. The hyper-Erlang and the Gaussian mixture distributions turned out to be the most versatile, providing the best fit in 35 and 10 cases, respectively. The node color in Figure 4 reflects the distribution family that fits the corresponding KPI the best.
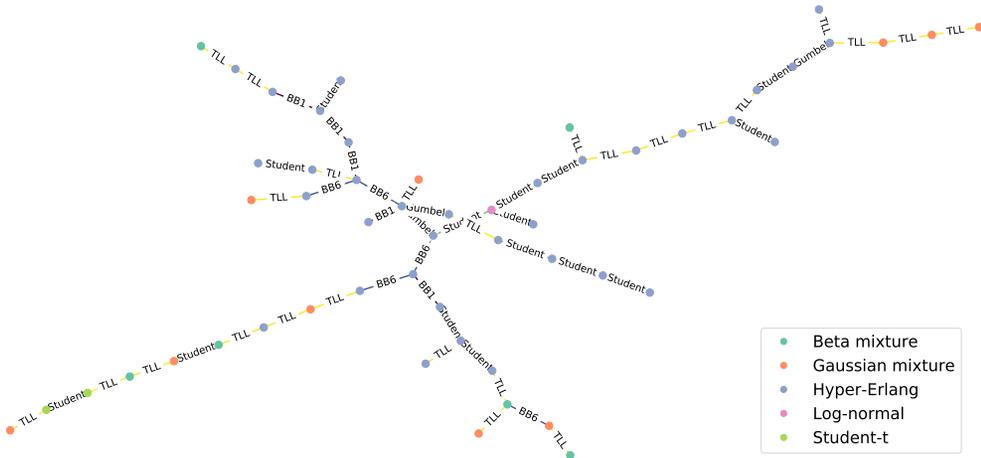
Fig. 4. Maximum information tree of the KPI pairs

A few examples for the marginal distribution fitting are shown in Figure 5 for some KPIs (these KPIs will be used later to demonstrate the copula and the joint distribution fitting). The results are excellent, the mixture distributions are flexible enough to capture not only the body, but also the tail of the empirical distributions. These examples represent the typical fitting accuracy, there are some better, and some worse results, too, but we managed to fit the majority of the marginal distributions with high accuracy.

### 6.3 Fitting bivariate copulas

To enable the separate fitting of the marginal distributions and the dependence structure, the two-dimensional subsets of the selected KPI pairs have been transformed to pseudo-observations. The resulting two-dimensional data has uniformly distributed marginals over [0, 1]. The two dimensional pdf obtained this way enables to infer the dependence structure of the two random variables.

To the transformed two-dimensional data several copula families are fitted (as detailed in Section 4.2), and the one leading to the highest likelihood is selected. Out of the candidate copula families the Gumbel copula turned out to be the best in 3 cases, the BB1 and the BB6 in 5 cases each, the Student copula in 15 cases, and the non-parametric TLL copula in 25 cases. The copula families associated with the pairs of variables linked by the edges of the maximum information tree are also available in Figure 4.

Figure 6 shows 3 examples for 3 copula families comparing the transformed observations with simulated samples of the fitted copulas. In the top row the BB1 copula, in the middle the Student copula, and in the bottom row the TLL copula provides the best fit. According to the examples in the Figure (that represent typical cases), the empirical copulas occurring in our dataset can be approximated with the considered copula families with high accuracy.

### 6.4 The bivariate joint distributions

Finally, we put together the bivariate copula and the marginal probability distributions for each KPI pair, reconstructing the bivariate joint probability distributions.
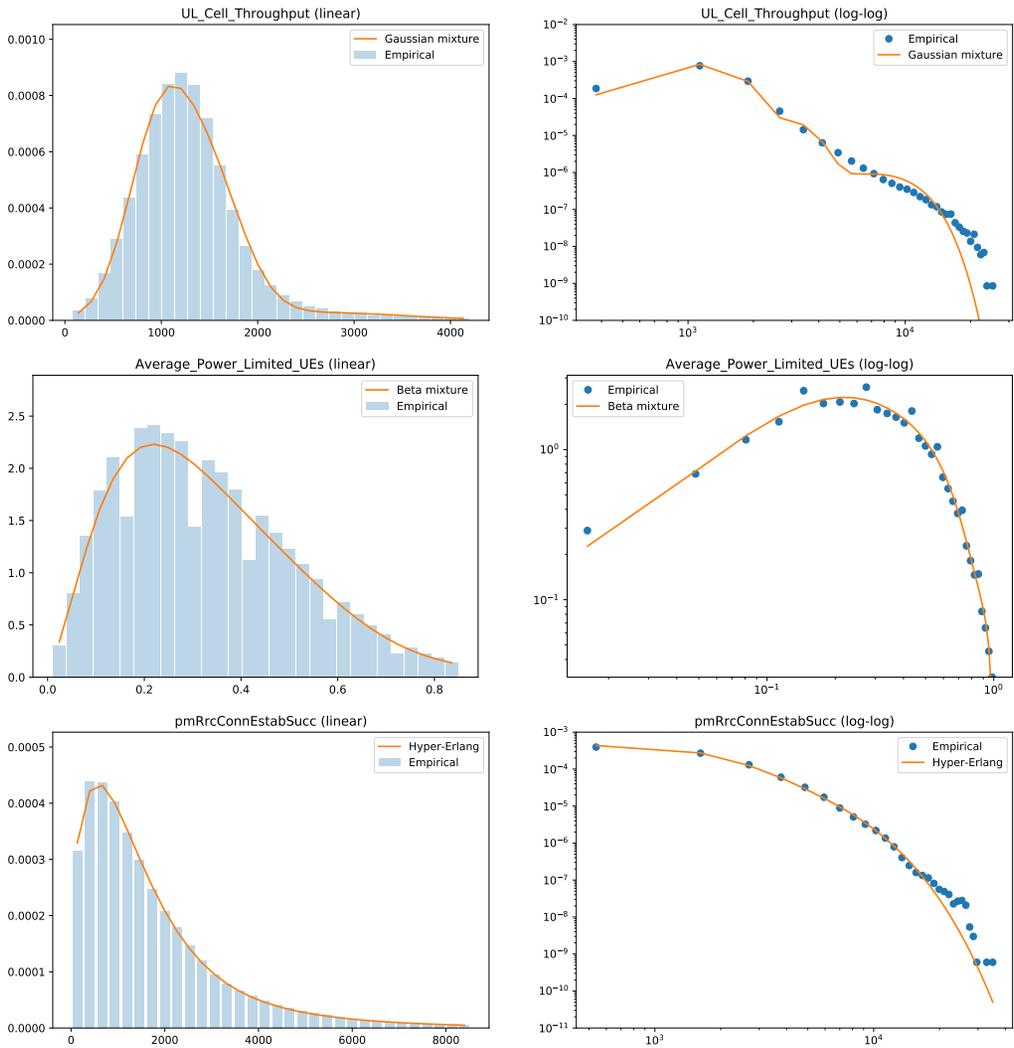
Fig. 5. The marginal distributions fitting KPIs UL_Cell_Throughput, Average_Power_Limited_UEs and pmRrcConnEstabSucc

The joint distributions corresponding to Figure 6 are depicted in Figure 7. In the figure, the original data points are shown in the left plot and random points simulated from the modelled joint distribution are given in the right side. The purpose of the contour lines is to make the comparison easier, they are generated from the probabilistic model and are the same for both the left and right hand side plots. The color of the points reflect the anomaly score (computed according to Section 5), light shades correspond to more rare (hence, more anomalous), dark shades to more typical observations. Observe that the joint distributions in these examples are very far from a bivariate Gaussian distribution, and are in fact difficult to fit with a mixture of bivariate Gaussians, too. However, based on the figures it is clear that the the copula-based modeling of the relevant variable pairs can be applied to this dataset successfully.

Fig. 6. Examples for copula fitting. Left: empirical copula. Right: simulated copula.

## 6.5 Anomaly tree

To visualize the current state of the network element (i.e., an LTE cell, in our case) we introduce *anomaly tree*s. The structure of the anomaly tree is given by the maximum information tree, and the weights of the edges are given by the anomaly scores calculated from the two dimensional marginal probability distributions. If the weights of the edges are visualized with different colors,

Fig. 7. Reconstructing the joint distributions. Left: original data points. Right: random points simulated from the model.

the anomaly trees provide a useful tool for the human operators to get an overview on the current state of the network element. Form engineering point of view the anomaly tree can be considered as a high level summary of the actual anomaly state of the monitored network element. The high weight edges identify the functional part(s) of the network element that are affected by the actual anomaly.

Fig. 8. Anomaly trees, example 1

Anomaly trees are illustrated in Figures 8 and 9. In these figures, the color indicates the severity of the anomaly. For example, the red color of a link means that the anomaly score is above 0.996. Orange shades mean lower anomaly score, and blue shades are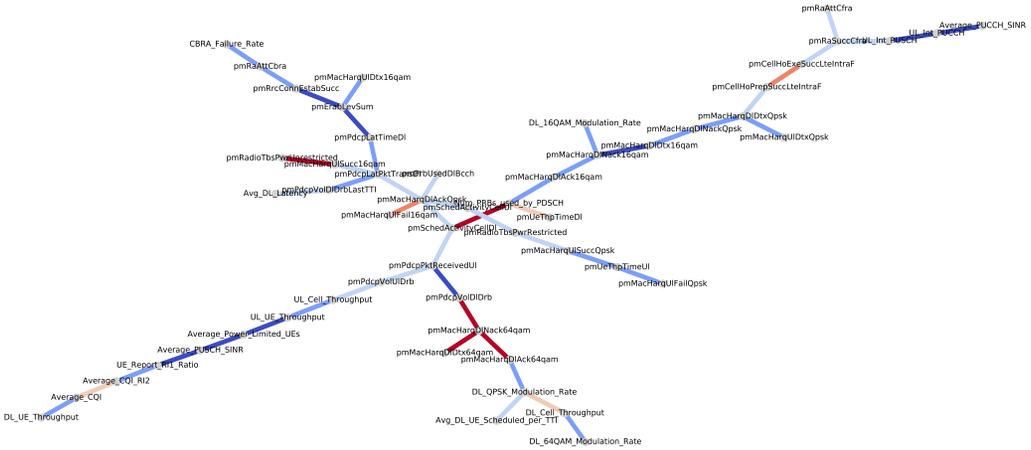 assigned to non-anomalous relations. Interpreting these anomaly trees requires deep domain specific knowledge, the detailed discussion of the possible reasons and the solutions is out of scope of this paper.

In Figure 8 we can see that the number of PRBs (Physical Resource Blocks) is extremely low, compared to the number of times a user equipment is selected for transmission (affecting `pmSchedActivityCellDl` - `Num_PRBs_used_by_PDSCH` relations). Given the transport block size (`pmRadioTbsPwrUnrestricted`), the high usage of the 16-qam modulation scheme (`pmMacHarq-UlSucc16qam`) and the large number of negative acknowledgements in case of the 64-qam modulation scheme (`pmMacHarqDlNack64qam`) implies that there could be something wrong with the radio channel, possibly with the antenna.

Figure 9 depicts an other kind of problem. In this case the channel quality is high enough, yet the user throughput is too low (according to the `Average_CQI` - `DL_UE_Throughput` relation). Considering the number of e-RABs (evolved radio access bearer), the number of RRC connection establishments is high (`pmErabLevSum` - `pmRrcConnEstabSucc`), and the signalling traffic is unusually high in upload direction (`UL_Int_PUCCH` - `UL_Int_PUSCH`), which means that the reason for the bad user experience is not the quality of the channel, but the sub-optimal setting of the RRC level parameters in the cell.

The manufacturer of the base station hardware has released guidelines for performance management and optimization, and also a troubleshooting guide, that the operators can use in order to resolve the problems detected by the presented anomaly detection algorithm.

## 6.6 Comparison with alternative methods

In this section we demonstrate how well the presented method performs in comparison with some standard and recent anomaly detection methods.

First we ignore the most appealing capability of the presented method, that is, the anomaly localization, and extract one overall anomaly score from the the whole anomaly tree. For an
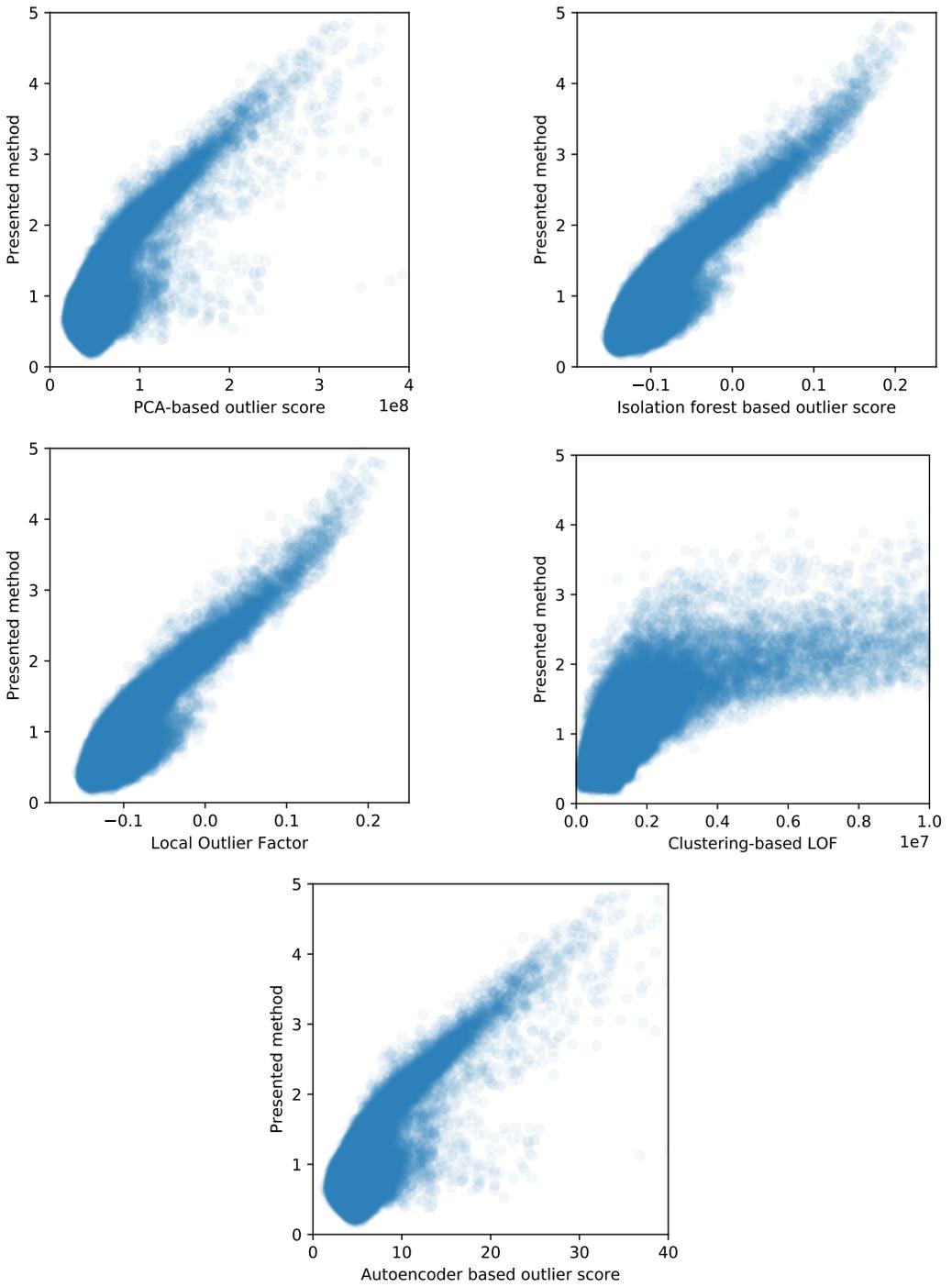
Fig. 9. Anomaly trees, example 2

observation **x** this overall anomaly score is obtained by

$$\tilde{\mathcal{A}}(\mathbf{x}) = \frac{1}{d-1} \sum_{e \in T} -\log(1 - \hat{\mathcal{A}}(\mathbf{x_e})), \tag{12}$$

where $\mathbf{x_e}$ is the two-dimensional sample containing only those two components of **x** that correspond to edge $e$ of the tree $T$. The role of the logarithm is to adjust the scale of the anomaly score such that only really rare observations will get high scores. Hence, after the logarithm re-scaling, we take the average anomaly scores of the anomaly tree.

Then, we compare this overall anomaly score to anomaly scores obtained by alternative methods. It is important to note that all the alternative methods we tested work on complete data only, and the LTE network data set contains only 66, 000 complete samples out of the 1, 500, 000 observations. Thus, in the rest of this section all comparisons are based on this very restricted, but complete data set.

To obtain anomaly scores from alternative procedures we have used the Python Toolbox for Scalable Outlier Detection (PyOD, [64]). Several procedures failed to provide a result in a reasonable time, or gave constant zero anomaly score for all samples (including the angle-based method [34]). For a comparison, Figure 10 shows the outlier scores of the presented methods against the scores of all the procedures that returned valid result, namely PCA-based outlier score [58], Isolation forest [40], Local Outlier Factor [10], Clustering-based LOF [27] and Autoencoder based outlier score [3]. For this study we have used the implementation available at https://github.com/yzhao062/pyod with the default parameters, as listed by Table 1.

From Figure 10 it can be seen that our method assigns high anomaly score to those observations that are found anomalous by alternative methods as well. Additionally, our method can provide information on the location of the anomaly and can operate with missing data as well. The fact that in a real data set missing data occurs frequently (possibly due to unreliable sensors) justifies the need for anomaly detection procedures that tolerate missing data.

Fig. 10. Comparison with alternative anomaly scoring methods

| Anomaly detection method | Parameters |
|---|---|
| PCA based method | use all components |
| Isolation forest | n_estimators=100 |
| Local Outlier Factor | n_neighbors=20 |
| Clustering-based LOF | n_clusters=8, $\alpha = 0.9$, $\beta = 5$ |
| Autoencoder based method | neurons=[64, 32, 32, 64], activation=relu, epochs=20 |

Table 1. Parameters of the methods involved in the comparison

## 7 EVALUATION OF THE METHOD ON A SYNTHETIC DATA SET

We have prepared an other numerical study to compare the presented copula based method to alternative anomaly detection methods published in the literature.

In this study we generate one million samples from a known distribution, and investigate how the anomaly scores correlate with the probability densities. For this purpose, we have defined a mixture of 5-dimensional normal distribution consisting of 4 components. The exact parameters of the distribution are provided in Appendix A. This choice may seem to be unfair, since our method is able to capture its marginals with high accuracy, but it turns out that this distribution has a very complicated correlation structure, that is equally difficult for all procedures involved in the comparison. As shown in Figure 11, a bi-variate joint distribution is multi-modal, our model can still capture its characteristic rather well.
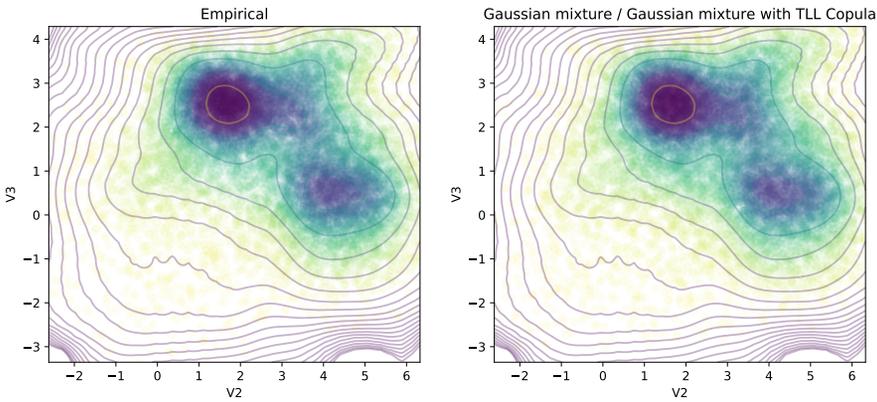


Fig. 11. One of the bi-variate joint distributions in the synthetic example. Left: original data points. Right: random points simulated from the model.

The procedures involved in the comparison are the same as in Section 6.6, namely the PCA based, the isolation forest based, the Local Outlier Factor based, the Clustering-based LOF and the Autoencoder based methods. Since the samples are generated from a know distribution, we can calculate the Kendall rank correlation coefficient between the anomaly scores returned by these methods and the values of the probability density function. The results are summarized in Table 2. The correlations are negative, since higher anomaly scores are assigned to observations having lower densities. According to the results, the presented copula based method performs the best for this particular data set.

We emphasize again, that this does not mean that the copula based method is always better than other methods, it just means that for data sets consisting of a high number of continuous

| Anomaly detection method | Kendall's tau |
|---|---|
| PCA based scoring | -0.337 |
| Isolation forest | -0.526 |
| Local Outlier Factor | -0.4194 |
| Clustering-based LOF | -0.5892 |
| Autoencoder based method | -0.3815 |
| **Presented method** | **-0.6275** |

Table 2. The correlation between the anomaly scores and the probability densities

variables it performs favorably. The alternative methods still have their strengths, when the data set is smaller, has categorical features and probability masses.

For completeness, we include the plots depicting the correlation between the copula based anomaly scores and the anomaly scores of the alternative methods, see Figure 12.

## 8 CONCLUSION AND DISCUSSION

In this paper, we presented a novel copula-based anomaly detection technique for high-dimensional data that does not only assign an anomaly score to the observations but it also localizes the reason of the anomaly. The proposed approach relies on the modeling of the multivariate probability distribution associated with the instances. Our procedure is able to handle large-scale, high-dimensional data since the intractable high-dimensional problem is broken into smaller tractable ones by using two dimensional projections (bivariate marginal distributions) of the joint probability distribution in such a way that it retains maximum information and reduces redundancy. Since rare events occur in the tails of the probability distribution, copulas were used to model the bivariate marginals. Another advantage of the copula approach that the univariate marginal probability distributions and the bivariate copulas can be fitted separately.

Besides an overall anomaly score, our approach also reports individual anomaly scores for the selected variable pairs that can be illustrated on an anomaly tree that enables the users to get an overview on the current state of the system, and to observe it evolving in time. We also illustrated our method on a real-world telecommunication problem. Our findings were reconfirmed by real network operators. Moreover, we compared the anomaly score of the proposed approach with other standard anomaly detection methods. We can conclude that the overall anomaly score returned by our approach is similar to the scores returned by other methods. On the other hand, our method also provides support to localize the anomaly, can be parallelized, can cope with a large number of high-dimensional observations and is also capable of handling missing data.

Although our proposed technique complements the existing approaches in several aspects, it also has its limitations. The most important limitation of the proposed approach is that it only works well with large scale continuous data. Categorical or discrete variables must be omitted since our method is not able to handle them (because the transformation to pseudo-observations does not result in a uniform distribution). This limitation also makes it troublesome to carry out a comprehensive comparison of the performance of our proposed approach with other alternative methods. Therefore, the method introduced in this paper supplements and does not substitute the existing techniques.

Fig. 12.  Comparison with alternative anomaly scoring methods on the synthetic data set

## REFERENCES

[1] Kjersti Aas, Claudia Czado, Arnoldo Frigessi, and Henrik Bakken. 2009. Pair-copula constructions of multiple dependence. *Insurance: Mathematics and economics* 44, 2 (2009), 182–198.

[2] Charu C Aggarwal. 2013. High-Dimensional Outlier Detection: The Subspace Method. In *Outlier Analysis*. Springer, 135–167.

[3] Charu C Aggarwal. 2015. Outlier analysis. In *Data mining*. Springer, 237–263.

[4] Fabrizio Angiulli and Clara Pizzuti. 2002. Fast outlier detection in high dimensional spaces. In *European Conference on Principles of Data Mining and Knowledge Discovery*. Springer, 15–27.

[5] Marcelo Bacher, Irad Ben-Gal, and Erez Shmueli. 2017. An Information Theory Subspace Analysis Approach with Application to Anomaly Detection Ensembles. (2017).

[6] Tim Bedford and Roger M Cooke. 2001. Probability density decomposition for conditionally dependent random variables modeled by vines. *Annals of Mathematics and Artificial intelligence* 32, 1-4 (2001), 245–268.

[7] Tim Bedford and Roger M Cooke. 2002. Vines: A new graphical model for dependent random variables. *Annals of Statistics* (2002), 1031–1068.

[8] Eike Brechmann. 2010. Truncated and simplified regular vines and their applications. (2010).

[9] Eike C Brechmann, Claudia Czado, and Kjersti Aas. 2012. Truncated regular vines in high dimensions with application to financial data. *Canadian Journal of Statistics* 40, 1 (2012), 68–85.

[10] Markus M Breunig, Hans-Peter Kriegel, Raymond T Ng, and Jörg Sander. 2000. LOF: identifying density-based local outliers. In *ACM sigmod record*, Vol. 29. ACM, 93–104.

[11] Guilherme O Campos, Arthur Zimek, Jörg Sander, Ricardo JGB Campello, Barbora Micenková, Erich Schubert, Ira Assent, and Michael E Houle. 2016. On the evaluation of unsupervised outlier detection: measures, datasets, and an empirical study. *Data Mining and Knowledge Discovery* 30, 4 (2016), 891–927.

[12] Varun Chandola, Arindam Banerjee, and Vipin Kumar. 2009. Anomaly detection: A survey. *ACM computing surveys (CSUR)* 41, 3 (2009), 15.

[13] Xiao-yun Chen and Yan-yan Zhan. 2008. Multi-scale anomaly detection algorithm based on infrequent pattern of time series. *J. Comput. Appl. Math.* 214, 1 (2008), 227–237.

[14] Umberto Cherubini, Elisa Luciano, and Walter Vecchiato. 2004. *Copula methods in finance*. John Wiley & Sons.

[15] C Chow and Cong Liu. 1968. Approximating discrete probability distributions with dependence trees. *IEEE Transactions on Information Theory* 14, 3 (1968), 462–467.

[16] Stephan Clémençon and Albert Thomas. 2017. Mass Volume Curves and Anomaly Ranking. *arXiv preprint arXiv:1705.01305* (2017).

[17] Xuan Hong Dang, Barbora Micenková, Ira Assent, and Raymond T Ng. 2013. Local outlier detection with interpretation. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 304–320.

[18] Rémi Domingues, Maurizio Filippone, Pietro Michiardi, and Jihane Zouaoui. 2018. A comparative evaluation of outlier detection algorithms: Experiments and analyses. *Pattern Recognition* 74 (2018), 406–421.

[19] John HJ Einmahl and David M Mason. 1992. Generalized quantile processes. *The Annals of Statistics* (1992), 1062–1078.

[20] Maria Erdmann. 2018. *Unsupervised Anomaly Detection in Sensor Data used for Predictive Maintenance*. Ph.D. Dissertation.

[21] Wentao Fan, Nizar Bouguila, and Djemel Ziou. 2011. Unsupervised anomaly intrusion detection via localized bayesian feature selection. In *2011 IEEE 11th International Conference on Data Mining*. IEEE, 1032–1037.

[22] Gabriel Frahm, Markus Junker, and Alexander Szimayer. 2003. Elliptical copulas: applicability and limitations. *Statistics & Probability Letters* 63, 3 (2003), 275–286.

[23] Javier Nunez Garcia, Zoltan Kutalik, Kwang-Hyun Cho, and Olaf Wolkenhauer. 2003. Level sets and minimum volume sets of probability density functions. *International Journal of Approximate Reasoning* 34, 1 (2003), 25–47.

[24] Supriya Garule and Sharmila M Shinde. 2015. Outliers Detection using Subspace Method: A Survey. *International Journal of Computer Applications* 112, 16 (2015).

[25] Gery Geenens, Arthur Charpentier, Davy Paindaveine, et al. 2017. Probit transformation for nonparametric kernel estimation of the copula density. *Bernoulli* 23, 3 (2017), 1848–1873.

[26] Markus Goldstein and Seiichi Uchida. 2016. A comparative evaluation of unsupervised anomaly detection algorithms for multivariate data. *PloS one* 11, 4 (2016), e0152173.

[27] Zengyou He, Xiaofei Xu, and Shengchun Deng. 2003. Discovering cluster-based local outliers. *Pattern Recognition Letters* 24, 9-10 (2003), 1641–1650.

[28] Harry Joe. 1997. *Multivariate models and multivariate dependence concepts.* CRC Press.

[29] JooSeuk Kim and Clayton D Scott. 2012. Robust kernel density estimation. *Journal of Machine Learning Research* 13, Sep (2012), 2529–2565.

[30] Edith Kovács and Tamás Szántai. 2016. Hypergraphs in the characterization of regular vine copula structures. *arXiv preprint arXiv:1604.02652* (2016).

[31] Edith Kovács and Tamás Szántai. 2017. On the connection between cherry-tree copulas and truncated R-vine copulas. *Kybernetika* 53, 3 (2017), 437–460.

[32] Hans-Peter Kriegel, Peer Kröger, Erich Schubert, and Arthur Zimek. 2009. Outlier detection in axis-parallel subspaces of high dimensional data. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining.* Springer, 831–838.

[33] Hans-Peter Kriegel, Peer Kröger, and Arthur Zimek. 2009. Clustering high-dimensional data: A survey on subspace clustering, pattern-based clustering, and correlation clustering. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 3, 1 (2009), 1.

[34] Hans-Peter Kriegel, Arthur Zimek, et al. 2008. Angle-based outlier detection in high-dimensional data. In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and data mining.* ACM, 444–452.

[35] Dorota Kurowicka and Roger M Cooke. 2006. *Uncertainty analysis with high dimensional dependence modelling.* John Wiley & Sons.

[36] Longin Jan Latecki, Aleksandar Lazarevic, and Dragoljub Pokrajac. 2007. Outlier detection with kernel density functions. In *International Workshop on Machine Learning and Data Mining in Pattern Recognition.* Springer, 61–75.

[37] Rikard Laxhammar, Goran Falkman, and Egils Sviestins. 2009. Anomaly detection in sea traffic-a comparison of the gaussian mixture model and the kernel density estimator. In *2009 12th International Conference on Information Fusion.* IEEE, 756–763.

[38] Aleksandar Lazarevic and Vipin Kumar. 2005. Feature bagging for outlier detection. In *Proceedings of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining.* ACM, 157–166.

[39] Kun-Lun Li, Hou-Kuan Huang, Sheng-Feng Tian, and Wei Xu. 2003. Improving one-class SVM for anomaly detection. In *Machine Learning and Cybernetics, 2003 International Conference on*, Vol. 5. IEEE, 3077–3081.

[40] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2008. Isolation forest. In *Data Mining, 2008. ICDM'08. Eighth IEEE International Conference on.* IEEE, 413–422.

[41] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2010. On detecting clustered anomalies using SCiForest. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases.* Springer, 274–290.

[42] Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. 2012. Isolation-based anomaly detection. *ACM Transactions on Knowledge Discovery from Data (TKDD)* 6, 1 (2012), 3.

[43] Ninghao Liu, Donghwa Shin, and Xia Hu. 2017. Contextual outlier interpretation. *arXiv preprint arXiv:1711.10589* (2017).

[44] Stephen Marsland, Jonathan Shapiro, and Ulrich Nehmzow. 2002. A self-organising network that grows when required. *Neural networks* 15, 8-9 (2002), 1041–1058.

[45] Emmanuel Muller, Ira Assent, Uwe Steinhausen, and Thomas Seidl. 2008. OutRank: ranking outliers in high dimensional data. In *2008 IEEE 24th international conference on data engineering workshop.* IEEE, 600–603.

[46] Roger B Nelson. 2007. Extremes of nonexchangeability. *Statistical Papers* 48, 2 (2007), 329–336.

[47] Hoang Vu Nguyen, Hock Hee Ang, and Vivekanand Gopalkrishnan. 2010. Mining outliers with ensemble of heterogeneous detectors on random subspaces. In *International Conference on Database Systems for Advanced Applications.* Springer, 368–383.

[48] Hoang Vu Nguyen and Vivekanand Gopalkrishnan. 2010. Feature extraction for outlier detection in high-dimensional spaces. In *Feature Selection in Data Mining.* 66–75.

[49] Hoang Vu Nguyen, Emmanuel Muller, and Klemens Bohm. 2013. 4s: Scalable subspace search scheme overcoming traditional apriori processing. In *Big Data, 2013 IEEE International Conference on.* IEEE, 359–367.

[50] Hoang Vu Nguyen, Emmanuel Müller, Jilles Vreeken, Fabian Keller, and Klemens Böhm. 2013. CMI: An information-theoretic contrast measure for enhancing subspace cluster and outlier detection. In *Proceedings of the 2013 SIAM International Conference on Data Mining.* SIAM, 198–206.

[51] Guansong Pang, Longbing Cao, Ling Chen, Defu Lian, and Huan Liu. 2018. Sparse modeling-based sequential ensemble learning for effective outlier detection in high-dimensional numeric data. In *Thirty-Second AAAI Conference on Artificial Intelligence.*

[52] Ninh Pham and Rasmus Pagh. 2012. A near-linear time approximation algorithm for angle-based outlier detection in high-dimensional data. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 877–885.

[53] Wolfgang Polonik. 1995. Measuring mass concentrations and estimating density contour clusters-an excess mass approach. *The Annals of Statistics* (1995), 855–881.

[54] Benjamin J Radford, Leonardo M Apolonio, Antonio J Trias, and Jim A Simpson. 2018. Network traffic anomaly detection using recurrent neural networks. *arXiv preprint arXiv:1803.10769* (2018).

[55] Ulf Schepsmeier, Jakob Stoeber, Eike Christian Brechmann, Benedikt Graeler, Thomas Nagler, Tobias Erhardt, Carlos Almeida, Aleksey Min, Claudia Czado, Mathias Hofmann, et al. 2018. Package VineCopula. (2018).

[56] Christopher Schröder and Sven Rahmann. 2017. A hybrid parameter estimation algorithm for beta mixtures and applications to methylation state classification. *Algorithms for Molecular Biology* 12, 1 (2017), 21.

[57] B Schweizer and A Sklar. 1969. Measures aleatoires de l'information. *CR Acad. Sci. Paris Ser. A* 269 (1969), 721–723.

[58] Mei-Ling Shyu, Shu-Ching Chen, Kanoksri Sarinnapakorn, and LiWu Chang. 2003. *A novel anomaly detection scheme based on principal component classifier*. Technical Report. Miami Univ Coral Gables FL Dept of Electrical and Computer Engineering.

[59] Lionel Tarassenko, Paul Hayton, Nicholas Cerneaz, and Michael Brady. 1995. Novelty detection for the identification of masses in mammograms. (1995).

[60] Axel Thummler, Peter Buchholz, and Miklós Telek. 2005. A novel approach for fitting probability distributions to real trace data with the EM algorithm. In *2005 International Conference on Dependable Systems and Networks (DSN'05)*. IEEE, 712–721.

[61] Nguyen Xuan Vinh, Jeffrey Chan, Simone Romano, James Bailey, Christopher Leckie, Kotagiri Ramamohanarao, and Jian Pei. 2016. Discovering outlying aspects in large datasets. *Data mining and knowledge discovery* 30, 6 (2016), 1520–1555.

[62] Xiaodan Xu, Huawen Liu, Li Li, and Minghai Yao. 2018. A Comparison of Outlier Detection Techniques for High-Dimensional Data. *International Journal of Computational Intelligence Systems* 11, 1 (2018), 652–662.

[63] Jifu Zhang, Sulan Zhang, Kai H Chang, and Xiao Qin. 2014. An outlier mining algorithm based on constrained concept lattice. *International Journal of Systems Science* 45, 5 (2014), 1170–1179.

[64] Yue Zhao, Zain Nasrullah, and Zheng Li. 2019. PyOD: A Python Toolbox for Scalable Outlier Detection. *arXiv preprint arXiv:1901.01588* (2019). https://arxiv.org/abs/1901.01588

[65] Arthur Zimek, Erich Schubert, and Hans-Peter Kriegel. 2012. A survey on unsupervised outlier detection in high-dimensional numerical data. *Statistical Analysis and Data Mining: The ASA Data Science Journal* 5, 5 (2012), 363–387.

## A  THE PARAMETERS OF THE DISTRIBUTION USED IN SECTION 7

To generate the syntetic data set for Section 7 we have defined a mixture of 5-dimensional normal distribution with 4 components, having density function

$$f(\mathbf{x}) = \sum_{i=1}^{4} p_i \frac{1}{\sqrt{(2\pi)^5 |\Sigma_i|}} e^{-(\mathbf{x}-\mu_i)^T \Sigma_i^{-1}(\mathbf{x}-\mu_i)/2}.$$

The mean values of the components are given by

$$\mu_1 = \begin{bmatrix} 2.9 & 4.3 & 0.6 & 4.0 & 4.4 \end{bmatrix}^T, \mu_2 = \begin{bmatrix} 4.0 & 1.5 & 2.7 & 0.7 & 4.5 \end{bmatrix}^T,$$
$$\mu_3 = \begin{bmatrix} 3.8 & 1.1 & 0.5 & 2.8 & 0.4 \end{bmatrix}^T, \mu_4 = \begin{bmatrix} 3.7 & 2.5 & 2.5 & 2.0 & 0.7 \end{bmatrix}^T,$$

and the mixing probabilities are $p_1 = 0.4, p_2 = 0.2, p_3 = 0.1, p_4 = 0.3$.

Obtaining the covariance matrices is a bit more involved. Our goal is to define the eigenvalues of the covariance matrices, and rotate them in a random way with unitary matrices. The eigenvalues of the covariance matrices $i = 1, \ldots, 4$, denoted by $\mathbf{v}_i$, are given by

$$\mathbf{v}_1 = \begin{bmatrix} 0.25 & 0.5 & 1 & 2 & 4 \end{bmatrix}, \mathbf{v}_2 = \begin{bmatrix} 0.15 & 0.25 & 0.35 & 0.5 & 1 \end{bmatrix},$$
$$\mathbf{v}_3 = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \end{bmatrix}, \quad \mathbf{v}_4 = \begin{bmatrix} 0.15 & 0.25 & 1 & 4 & 5 \end{bmatrix}.$$

Before providing the covariance matrices, let us introduce the following auxiliary matrices:

$$
\mathbf{W}_1 = \begin{bmatrix}
0.6 & 0.1 & 0.9 & 0.9 & 0.4 \\
0.8 & 0.7 & 0.2 & 0.8 & 0.6 \\
0.8 & 0.1 & 0.5 & 0.3 & 0.4 \\
0.8 & 0.6 & 0.7 & 0.7 & 0.8 \\
0.6 & 0.6 & 0.8 & 0.7 & 0.7
\end{bmatrix}, \mathbf{W}_2 = \begin{bmatrix}
0.2 & 0.4 & 0.5 & 0.6 & 0.1 \\
0.8 & 0.7 & 0.5 & 0.0 & 0.8 \\
0.9 & 0.2 & 0.2 & 0.8 & 0.8 \\
0.3 & 0.2 & 0.5 & 0.3 & 0.7 \\
0.8 & 0.8 & 0.8 & 0.3 & 0.5
\end{bmatrix},
$$

$$
\mathbf{W}_3 = \begin{bmatrix}
0.1 & 0.3 & 0.0 & 0.9 & 0.8 \\
0.5 & 0.7 & 0.5 & 0.6 & 0.1 \\
0.0 & 0.3 & 0.2 & 0.3 & 0.2 \\
0.5 & 0.7 & 0.8 & 0.5 & 0.6 \\
0.3 & 0.1 & 0.9 & 0.6 & 0.1
\end{bmatrix}, \mathbf{W}_4 = \begin{bmatrix}
0.4 & 0.4 & 0.3 & 0.2 & 0.5 \\
0.7 & 0.4 & 0.3 & 0.4 & 0.6 \\
0.0 & 0.8 & 0.3 & 0.9 & 0.8 \\
0.9 & 0.9 & 0.1 & 0.8 & 0.1 \\
0.1 & 0.3 & 0.3 & 0.9 & 0.1
\end{bmatrix}.
$$

Applying QR decomposition on matrices $\mathbf{W}_i, i = 1, \dots, 4$ gives orthogonal matrices $\mathbf{Q}_i, i = 1, \dots, 4$, from which the covariance matrices are obtained by $\mathbf{C}_i = \mathbf{Q}_i \operatorname{diag}(\mathbf{v}_i) \mathbf{Q}_i^T$.