

Automated Identification of Network Anomalies and Their Causes with Interpretable Machine Learning: the CIAN Methodology and TTrees Implementation¹

Mohamed Moulay^a, Rafael Garcia Leiva^b, Pablo J. Rojo Maroni^c, Fernando Diez^d, Vincenzo Mancuso^e, Antonio Fernández Anta^e

^aUniversity Carlos III of Madrid, Spain

^bVodafone, Madrid, Spain

^cNokia Cloud and Networks Services, Madrid

^dUniversidad Politécnica de Madrid, Spain

^eIMDEA Networks Institute, Madrid, Spain

Abstract

Leveraging machine learning (ML) for the detection of network problems dates back to handling call-dropping issues in telephony. However, troubleshooting cellular networks is still a manual task, assigned to experts who monitor the network around the clock. To help in this task we present CIAN (from Causality Inference of Anomalies in Networks), a practical and interpretable ML methodology, which we implement in the form of a software tool named TTrees (from Troubleshooting Trees). We have designed CIAN to automate the identification of the causes of performance anomalies in cellular networks. Our methodology is unsupervised and combines multiple ML algorithms (e.g., decision trees and clustering) and Kolmogorov complexity-inspired data analysis tools that we have developed for this work. CIAN can be used with small volumes of data and is quick at training.

Our experiments use diverse data sets obtained from measurements in operational commercial mobile networks. They show that the TTrees implementation of CIAN can automatically identify and accurately classify network anomalies—e.g., cases for which a network low performance is not apparently justified by operational conditions—training with just a few hundreds of data samples. The resulting information hence enables precise troubleshooting actions. In particular, we showcase how TTrees can be flexibly used to monitor the performance of TCP and QUIC protocols when they are adopted to serve mobile users.

Keywords: Troubleshooting, Anomaly detection, Feature selection, Interpretable machine learning

1. Introduction

Cellular networks have been through an exceptional evolution in recent years. With 4G, and even more with 5G, network services have gained a significant degree of intelligence, involving intensive access to both data communication and computing resources. With the evolution of cellular networks, there has also come an increase in structural complexity and heterogeneity of services, requiring constant monitoring of the communication system. This landscape will only become more complex with the evolution of 5G and the emergence of 6G [2]. Most importantly, new systems will require intelligent and automatic troubleshooting tools, which is the focus of our work.

A line of work on monitoring and troubleshooting of cellular networks is the development of self-healing networks. Self-healing networks are responsible for detecting, identifying, and making decisions on recovery actions [3]. Multiple proposals exist for making fault detection and self-healing systems practical in mobile networks [4]. However, current self-healing troubleshooting proposals lack flexibility and do not scale well. As an alternative, there are newly-defined approaches based on ML, which use deep learning and neural networks as black boxes [5]. Unfortunately, these approaches lack interpretability, so it is not possible to understand the cause of a detected network

¹A preliminary version of this work appeared in the proceedings of IEEE WoWMoM 2021 [1]

14 problem, and manual intervention is required for classification after the detection. It is then possible to resort to Ex-
15 plainable AI, or XAI, which deals with the problem of how human users could understand AI's cognition, and decide
16 if an AI based model can be trusted or not [6]. However, XAI does not help interpret AI decisions and conclusions *per*
17 *se*, although its transparency of operation and interpretability of decisions allows humans to identify which features
18 of the observed system led to specific algorithm's outputs on a case-by-case basis [7]. Multiple methods have been
19 proposed to address the complex issue of ML interpretability, from determining which features contribute the most
20 to a neural network's output, to the development of targeted models that explain individual predictions [8]. However,
21 other studies point out that there is a potential need for applying an interdisciplinary joint efforts in order to correctly
22 apply and leverage explainable/interpretable ML algorithms [9].

23 In summary, problems like troubleshooting remain a substantially manual procedure, in which highly skilled
24 experts analyze alarms and statistics of performance indicators regularly. These experts need to interpret the outcome
25 of statistical models and ML routines in order to detect and diagnose the cause of problems and anomalies in the
26 network. In contrast, unskilled engineers can not even detect problems effectively [10].

27 In this paper we focus on detecting and determining the causes of network anomalies automatically. We do
28 not focus on identifying problems for which behavioral models exists or can be built. Instead, we only focus on
29 anomalies that models cannot explain. Specifically, we define an anomaly as a situation that should not happen, given
30 the capacity of the network and the current network traffic. Anomalies are not the same as problems. For instance,
31 while an overloaded base station located in a football stadium full of people yields a network problem, that is not
32 an anomaly, because there is a clear reason for the base station to provide low per-user throughput. On the contrary,
33 an idle base station when the stadium is full is not a network problem (no network alarm is raised), but it is an
34 anomaly, since the base station should be busy. This situation should raise an alarm and call for examination in order
35 to seek for the causes of this anomaly, which might go beyond the network. We are interested in the identification
36 of anomalies because they complement our tools for the identification of problems across the components used to
37 establish a communication between users and servers in a network. Moreover, anomalies are, hopefully, easier to
38 solve than network problems, since they deal with an abnormal operation of the network, not with its limitations.

39 We go beyond mere anomaly detection and provide an interpretable methodology that can identify the causes of
40 anomalies within the observable features of the system. Our proposed anomaly detection process is applied to data
41 collected in real data from drive test of networks. Specifically, to benchmark QoS in mobile networks, continuous
42 drive tests with end-to-end test scenarios are performed every day internally by the network operators, and externally
43 by third-parties or government regulators. Each of these test campaigns can have up to tens of thousands of individual
44 test cases, from which specific KPIs are calculated. Drive tests are normally executed with off-the-shelf testing equip-
45 ment (NEMO, TEMS, Swissqual, etc.), capable of running predefined sequences of tests and collecting relevant low
46 level radio and traffic information, as well as application performance statistics. Those test harvest many operational
47 parameters (features) and key performance indicators (KPIs) that allow to monitor network health, so as to evaluate
48 the Quality of Service (QoS) that the network is providing to customers. Features and KPIs identify the operational
49 context of the network and measure the performance of radio, TCP, routing, etc. Hence, we use features and KPIs to
50 detect and classify anomalies, and to trigger troubleshooting actions.

51 To perform anomaly detection and classification of their causes, we resort to simple and unsupervised ML tools to
52 process the data collected from the network. Using these ML tools, we build an interpretable and robust methodology,
53 named CIAN, which we have implemented as a software named TTrees and tested with real data. In addition to
54 interpretability, a key advantage of CIAN and TTrees is that it does not require access to large volumes of training
55 data.

56 1.1. Original Contribution of the Work

57 We propose an unsupervised methodology for the detection and classification of network performance anomalies,
58 which we call Causality Inference of Anomalies in Networks (CIAN). In addition, we provide a software implemen-
59 tation of the methodology named Troubleshooting Trees (TTrees). Hence, with this work, we join the ML research
60 stream in networking, while focusing on the troubleshooting of possible network issues even with small volumes of
61 measurement data. However, differently from existing ML proposals, we detect anomalies by simply identifying the
62 scenarios that ML algorithms *cannot* learn/explain. For instance, a network that shows low throughput is not anoma-
63 lous if its low performance can be explained, e.g., by detecting the presence of a bad radio link. For this reason,
64 a novel aspect of our work is the use of simple interpretable ML algorithms, moving away from the current trend

65 of high-accuracy ML algorithms (e.g., deep learning) that do not allow interpretation (and hence understanding) of
66 their outcome. Specifically, we use decision trees in our methodology since they are transparent to inspection, hence
67 interpretable [11]. In fact, having low accuracy ML models is not a problem to our system, since we consider as
68 interesting (anomalous) the scenarios that are misclassified by the ML algorithms, and we do not want to miss them
69 by overfitting. Then, the anomalous scenarios are classified by the potential causes of their behavior. Understanding
70 and classifying these anomalous scenarios allow alerting the appropriate department to take corrective actions.

71 In order to evaluate the methodology, we use real operational network data. Two of the datasets used here were
72 collected for cellular service auditing purposes by Nokia in various European countries. They include thousands of
73 features (i.e., operational parameter reports and statistics) and a few KPIs used to mark the quality of the network
74 in each test (i.e., throughput, connection establishment time, download time, etc.). However, the available datasets
75 are also quite limited in terms of number of data samples, e.g., they only report hundreds or a few thousands of
76 experiments each. Thus, these datasets are not suitable for commonly adopted deep neural networks and ML methods
77 requiring complex training. The other datasets have been obtained in measurement campaigns conducted by us with
78 MONROE, an open-access platform for multi-homed experimentation with commercial mobile broadband networks
79 across Europe [12]. With MONROE, we have access to larger datasets, although with a reduced number of features
80 with respect to Nokia’s datasets. Moreover, while Nokia’s datasets report mostly TCP traffic, with MONROE we were
81 able to experiment with TCP and QUIC. As a result, in this article we show how CIAN can be flexibly used to study
82 heterogeneous and diverse datasets to identify the root causes of anomalies without human intervention.

83 The experts in detection and classification of network issues from the research team have been instrumental in the
84 validation of the methodology and the TTrees system. They have manually inspected the outcome of the unsuper-
85 vised process, verifying that the scenarios that were classified as anomalous did in fact show strange combinations
86 of observed features and KPIs, and that their unsupervised classification into *network aspects* (i.e., combination of
87 observed operational features) that identify the causes of the anomalies was consistent. Additionally, the experts
88 guided the development of a supervised implementation of the methodology, STrees, in which features were manually
89 aggregated into a few relevant network aspects based on the experts’ knowledge and experience, like it happens in
90 currently implemented troubleshooting protocols. The comparison of the outcomes of unsupervised and supervised
91 methods with the same data has also been used to validate the methodology.

92 As mentioned, we provide a Python open-source implementation of our methodology based on the Scikit-learn
93 library [13] and a new library that we have implemented.²

94 In this extended version of the paper, we carried out an in-depth evaluation of the methodology and its implemen-
95 tation through the execution of experiments in cellular networks. More specifically, we expanded the related work by
96 introducing additional references on causality and explainable AI with the current state of the art, introduced and ana-
97 lyzed two datasets with imbalanced anomalies and tested an alternative implementation of CIAN with a new clustering
98 algorithm. We used this alternative implementation to draw performance comparisons that assess the shortcomings
99 of TTrees and propose possible solutions to the latter.

100 1.2. Organization of the Paper

101 The rest of the paper is organized as follows. Section 2 discusses the related work. Section 3 overviews CIAN,
102 the unsupervised methodology proposed. Section 4 explains in detail how this methodology has been implemented as
103 the TTrees system. Section 5 presents and analyzes the results TTrees achieves when it is applied to real datasets of
104 various kind. Finally, we summarize and conclude the paper in Section 6.

105 2. Related Work

106 2.1. Causality

107 Causality refers to the relationship existing between two types of events, cause and effect, where the occurrence of
108 the latter is a consequence of the occurrence of the first. Causality must not be confused with another common term
109 belonging to the statistics lingo, which is correlation. While the latter refers to the relationship between two variables

²<https://github.com/Mohmoulay/WoWMoM2021>

110 based on their covariance (if a variable changes, the other does so too), causality establishes a clear dependency and a
111 mapping in the way these changes take place. Causality always implies correlation but not the other way around [14].

112 Traditional methods for the discovery of causal relationships include the conduction of randomized trials with
113 the objective of removing confounders, that is, variables that may deceive into drawing non-existent associations
114 between different events. Due to the high cost of implementing these experiments, many researches have switched to
115 the discovery of causal relationships through the exploratory analysis of observational data. This is known as causal
116 inference [15]. As a consequence, and due to the exponential increase in the amount of data that can be collected
117 and analyzed by systems nowadays, machine learning has been put at the forefront as a promising pathway to the
118 automatization of causal inference and its execution in a computationally feasible amount of time. However, the
119 introduction of machine learning models can lead to the so-called “black box effect,” as the convoluted nature of
120 these algorithms can obscure the understanding of the reasons behind the cause-effect connections and predictions
121 discovered [16].

122 Advances have been done theoretically to favour the explanation of these causal relationships in the field of
123 machine learning. Judea Pearl proposed a three-level hierarchy to classify causal information according to the type
124 of queries each family can give answer to [17]. Firstly, the Association level, at the bottom of the causal hierarchy
125 and syntactically expressed through conditional probability sentences, deals with information that allows to infer
126 data associations extracted from the pure observation of data and the application of standard probabilities. Secondly,
127 the Intervention level, regards questions that not only imply the observation of data but also the alteration of the
128 perceived reality. The probabilistic expressions at this level can be estimated both through experiments and analysis
129 using causal Bayesian networks. Finally, at the top of the hierarchy we have the Counterfactual level, which deals
130 with retrospective reasoning and the hypothesis of scenarios provided environmental conditions had been different or
131 altered prior to observation.

132 Pearl’s three-tier hierarchy sheds light on why most machine learning systems, which are only capable of extract-
133 ing data associations and thus fit in the bottom of the causal hierarchy, are unable to reason about newly unobserved
134 data and provide causal explanations. Thus, this raises the question whether artificial intelligence can be enriched to
135 provide new layers of causal inference. The main pathway towards a second generation of machine learning models
136 capable of achieving this level of cognition involves the development of Explainable AI (XAI).

137 2.2. *Explainable AI*

138 Using Machine learning to detect anomalies has been around for a while. Currently, the most powerful algorithms
139 are based on Neural Networks, which show high accuracy but have little interpretability. In reality, interpretability
140 mainly refers to the intuition behind the outputs of a model; the more interpretable a machine learning system is, the
141 easier it is to identify causality within the system’s inputs and outputs. Thus, XAI has gained much momentum in
142 the past few years to help explain any black-box Model. The most complete techniques to achieve that are the local
143 interpretable model-agnostic explanations (LIME) and Shapley Additive explanations (SHAP), based on the current
144 state of the art.

145 The LIME [18] method is one of the most popular interpretability methods for black-box models. Following a sim-
146 ple yet powerful approach, LIME can generate interpretations for single prediction scores produced by any classifier.
147 For any given instance and its corresponding prediction, simulated randomly-sampled data around the neighborhood
148 of the input instance for which the prediction was produced are generated. Subsequently, new predictions are made
149 for the generated instances and weighted by their proximity to the input instance while using the model in question.
150 Lastly, a simple, interpretable model, such as a decision tree, is trained on this newly-created dataset of perturbed in-
151 stances. By interpreting this local model, the initial black-box model is consequently interpreted. Although LIME is
152 powerful and straightforward, it has its drawbacks. In 2020, the first theoretical analysis of LIME [18] was published,
153 validating the significance and meaningfulness of LIME and proving that poor parameters choices could lead LIME
154 to miss out on essential features.

155 SHAP [19] is a game-theory-inspired method that attempts to enhance interpretability by computing the essential
156 values for each feature for individual predictions. Firstly, the authors define the class of additive feature attribution
157 methods, which unifies six current methods, including LIME [18], DeepLIFT [20], and Layer-Wise Relevance Propa-
158 gation [21], which all use the same explanation model. Subsequently, they propose SHAP values as a suitable feature
159 importance measure that maintains three desirable properties: local accuracy, missingness, and consistency. Finally,

160 they present several different methods for SHAP value estimation and provide experiments demonstrating the supe-
161 riority of these values in terms of differentiating among the different output classes and better aligning with human
162 intuition than many other existing methods.

163 These techniques are generic and help interpret models deployed in multiple disciplines, such as healthcare,
164 medicine, retail, and banking. A key factor concerning these techniques is that they have to be used on pre-trained
165 models. We differ because we are trying to deploy techniques that are easy to interpret throughout the entire process,
166 from data processing to problem classification.

167 2.3. Anomaly detection

168 Early works on fault detection suggested the use of time series *regression* methods and *Bayesian networks*. For
169 instance, Barco *et al.* [22] produced an automated tool for troubleshooting mobile communication networks back in
170 the days of 2G, relying on Bayesian networks to detect call drops. Khanafer *et al.* [23] followed up on proposing
171 a method based on Bayesian networks to detect faults in UMTS systems, in which they apply different algorithms
172 to discrete KPIs. Other works, such as [24], rely on a scoring-based system, in which the authors build the fault
173 detection subsystem around labeled fault cases. These cases were previously identified by *experts*, using a scoring
174 system to determine how well a specific case matches each diagnosis target. The work presented in [25] is based
175 on a supervised genetic fuzzy algorithm that learns a fuzzy rule base and, as such, relies on the existence of labeled
176 training sets. Indeed, most of the techniques proposed in the literature focus on using supervised machine learning
177 algorithms [23, 24, 25]. In this paper we show that it is convenient to use unsupervised techniques to unveil hidden
178 information in the input data, without restricting a priori the possible outcomes.

179 Other works make use of advanced mathematical and statistical tools. For instance, Ciocarlie *et al.* [26] address
180 the problem of checking the effect of network changes via monitoring the state of the network, and determining if the
181 changes resulted in degradation. Their fault detection mechanism uses *Markov logic networks* to create probabilistic
182 rules that distinguish between different causes of problems. A framework for network monitoring and fault detection
183 is introduced in [27], using *principal component analysis* (PCA) for dimension reduction, and kernel-based semi-
184 supervised fuzzy clustering with an adaptive kernel parameter. To evaluate the algorithms, they use data generated by
185 means of an LTE system-level simulator. The authors claim that this framework proactively detects network anomalies
186 associated with various fault classes. These methods lack the flexibility of ML-based ones and, differently from our
187 proposal, cannot be fully automated for a generic network context.

188 Recently, researchers are exploring the potentials of AI/ML for predicting and timely obviating network perfor-
189 mance issues. For instance, Terra *et al.* [28] analyze how to apply existing XAI methods to identify the root causes of
190 service level agreement violations in a sliced network, in a 5G context. Tang *et al.* [29] use instead ML to identify the
191 root causes of exceptions in packet-over-optical networks, thus achieving what they define customized performance
192 monitoring and troubleshooting. Furthermore, Wei *et al.* [30] point out that the class of AI that goes under the label of
193 intent-based networking techniques can be used for troubleshooting of network services. For instance, available prod-
194 ucts like Cisco ACI and Spruce Network DeepFlow already offer network operators the capability to monitor and rise
195 alarms and trigger troubleshooting actions in 5G and beyond 5G networks. However, they do not offer explainable nor
196 interpretable methods and basically help in the self-optimization of a system using intent-based networking, which is
197 limited to the scope of SDN/NFV network functionalities, while we will see in the reminder of this article that con-
198 textual information is key to identify anomalies rather than generic performance issues, and their causes. Moreover,
199 existing AI/ML based approaches to troubleshooting are highly customized (e.g., for 5G functionalities or for optical
200 networks) and require the collection of specific network features. Differently from our approach, these limitations
201 make existing approaches unsuitable to be flexibly applied to new protocols and services.

202 Finally, it is worth mentioning that most of the existing proposals have been evaluated only through simulators, and
203 require large datasets. They help to detect network issues, but do not contribute to interpreting the network behavior.
204 By comparison, in our work, we use not-necessarily-abundant data collected in real operational networks and propose
205 a fully automated ML-based methodology that leads to a straightforward interpretation of network behaviors. This
206 includes identifying not only the occurrence of problems but also their root causes.

207 **3. Overview of CIAN**

208 In this section, we briefly describe the framework, and the steps into which the automated unsupervised method-
 209 ological process we propose is divided, without entering in the details of design and implementation, which will be
 210 instead the object of Section 4. We begin by defining some basic notations and concepts that will be used to describe
 the methodology. The basic notation used is summarized in Table 1.

Table 1: Basic notation.

Symbol	Description
M	Model
D	Dataset
$L(M)$	Model length
$L(D M)$	Length of the dataset D given the model M
\mathbf{x}	Collection of data samples
\mathbf{x}_i	A specific data sample
x_{ik}	A specific feature of a data sample
\mathbf{y}	Set of target KPIs
m	Number of classes resulting from KPI discretization
C_1	Initial classifier (Knowledge Tree)
A	Set of anomalous scenarios
\mathbf{x}^A	Set of features in anomalous scenarios
\mathbf{y}^A	Set of target KPIs in anomalous scenarios
R	Set of features that are highly informative about the target KPI
r_a	Subset of relevant features for problem classification
α	Number of subsets selected for problem classification
C_2	Second classifier (Aspect Classification Tree)

211

212 **3.1. The Core Idea for Detecting Anomalies**

213 Our methodology has been developed to find and interpret anomalies in network communications, in an automa-
 214 tizable way. The core idea behind our approach has been inspired by the concepts of *Kolmogorov complexity*, and in
 215 particular, by the *minimum description length* (MDL) and *minimum message length* principles [31, 32]. According to
 216 these concepts, learning from data is equivalent to find regularities, or equivalently, compress data, i.e., describe data
 217 in the shortest possible way without losing information. In particular, the MDL principle claims that the best model
 218 M for a dataset D is the one that minimizes the following expression.

$$L(M) + L(D|M), \tag{1}$$

219 where $L(M)$ is the length of a model M encoded as a string of symbols, and $L(D|M)$ is the length of the dataset
 220 given the model. Intuitively, this latter term can be seen as the length of the list of errors made by the model M
 221 when predicting KPIs from features with a specific dataset D , encoded using an optimal length code. In this sense,
 222 we divide our dataset into two disjoint groups: the compressible part of D that can be described through M , and
 223 the incompressible part that includes all the points in D that are wrongly predicted (e.g., misclassified) by M . In a
 224 network data context, the compressible part of D is therefore what the model M identifies as regularities in the network
 225 described with a given set of features and KPIs, while the incompressible part corresponds to unexpected behaviors for
 226 which the KPIs cannot be predicted based on the observed features. In general, the more complex we allow the model
 227 to be, the more regularities we will be able to find, but also the longer its representation will be; thus, the appropriate
 228 balance between model complexity and model accuracy has to be found, so as to minimize expression (1). In this
 229 paper we address this trade-off between model complexity and model accuracy by means of applying an incremental
 230 procedure over model hyper-parameters to find the most complex possible model that does not overfit the dataset. The
 231 particular families of models used are selected for their accuracy and interpretability, i.e., they produce models that
 232 identify the compressible part of D based on interpretable rules. From here it follows the claim that the incompressible
 233 part of D cannot be interpreted, and hence it contains *anomalies*.

234 3.2. Input

235 The input to the process is a collection \mathbf{x} of data samples, obtained from n network operation scenarios. Each
236 data sample \mathbf{x}_i includes k features (x_{i1}, \dots, x_{ik}) , which are operational characteristics of the network and the context
237 in which measurements are taken, in scenario $i \in [1, n]$. In addition to these features, there is a distinguished set of
238 ℓ target KPIs \mathbf{y} which drives the search for anomalies. Hence, the dataset is formed by \mathbf{x} and \mathbf{y} , and each scenario
239 $i \in [1, n]$ is a row of the dataset containing features (x_{i1}, \dots, x_{ik}) and KPIs $(y_{i1}, \dots, y_{i\ell})$. Note that realistic datasets can
240 contain incomplete data entries. Therefore, before processing further a dataset, we will possibly filter out malformed
241 or incomplete rows, and all rows that do not contain values for the target KPIs.

242 3.3. Discretization

243 The first step of our methodology consists in grouping data samples into classes, based on the value of the target
244 KPIs \mathbf{y} . The reason for having this step is because, in general, the target KPI values belong to an infinite domain.³
245 Note that our methodology is oblivious to the process that collects data samples, as we simply take a dataset and work
246 with all its valid entries. For discretization, it is desirable that the number of classes m used is manageable but not too
247 small, for expressiveness. However, we do not need vast amounts of data, as our methodology works with as few as
248 tens of samples per discretization class, so as to have a bare minimum level of statistical relevance. Additionally, it is
249 desirable that data samples are reasonably balanced across classes, although this may not be perfectly achievable in
250 all cases. Finally, in order to provide semantics to the discretization, every vector of KPIs $(y_{i1}, \dots, y_{i\ell})$ is assigned a
251 certain QoS, which can be quantified. Then, the discretization places in the the same class scenarios whose QoS are
252 similar, and the m classes can be ranked by an expert based on their goodness (e.g., as “very bad”, “bad”, “medium”,
253 “good”, “very good”, etc.).

254 3.4. Selection of Anomalous Scenarios

255 Once the data samples have been assigned to the classes based on the value of the target KPIs, we use them to
256 train a classifier C_1 . This classifier will use the features (x_{i1}, \dots, x_{ik}) to determine whether, in scenario i , the KPIs
257 $(y_{i1}, \dots, y_{i\ell})$ seem to belong to class $j \in [1, m]$. The objective is to build a classifier that is able to correctly classify
258 as many data samples as possible without overfitting. It is also desirable that the decisions of the classifier can be
259 interpreted.

260 Applying this classifier C_1 to all the data samples will hence incorrectly classify a subset of them, out of which we
261 consider the subset $A \subset [1, n]$ for which the classifier predicts classes higher than the observed ones. The scenarios
262 in set A are the ones that we consider interesting, or *anomalous*: scenarios whose target KPIs is overestimated with a
263 properly trained classifier C_1 .⁴ The interpretability of the classifier helps to determine why the data samples in A are
264 anomalous (with respect to those properly classified). This means identifying which features are relevant and how to
265 differentiate A from the rest.

266 Observe that the scenarios in A are not necessarily those in which any of the target KPIs show low performance.
267 For instance, if \mathbf{y} consists of a single KPI (e.g., the average throughput observed), a scenario i with low throughput
268 y_i may not be anomalous if this is due to a low radio coverage (i.e., when the user is at the cell edge). This scenario
269 should be correctly classified by an interpretable classifier C_1 as bad, and a visual inspection of the radio features in \mathbf{x}_i
270 should reveal the issue. The anomalous scenarios that are of interest to us are those that cannot be (directly) explained.
271 For instance, a scenario i with low throughput y_i in which the features \mathbf{x}_i are all good.

272 Hence, from this point on, our target will be to identify automatically which types of anomalies the scenarios in A
273 show. This could later be used to report the issue (complemented with the data samples) to the appropriate department,
274 that may take corrective action.

³We explored options without discretization, but they were unsatisfactory, because they led to complex approaches with additional hyperparameters.

⁴In some practical cases, it may be interesting to consider in the set A also the scenarios for which the classification assigns a worse class (in QoS terms) than their actual class. Anomalies that have a positive impact on QoS are not harmful at all for the user and could be treated less urgently than anomalies with negative impact on QoS, so we ignore them in this work. Still, they would deserve further studies because they could reveal the existence, and possibly explain the cause, of new interesting effects to be considered for design and tuning of network systems and services. The methodology described in this article can be straightforwardly extended to analyze also this kind of anomalies.

275 3.5. Selection of the Most Relevant Features

276 Let \mathbf{x}^A and \mathbf{y}^A be the set of features and corresponding KPIs in anomalous scenarios, taken as \mathbf{x}_i and \mathbf{y}_i , for $i \in A$.
277 These are the scenarios that we want to explore further. A natural approach to attempt for understanding what makes
278 these scenarios anomalous would be using another classifier only for them. While this indeed separates these scenarios
279 into different classes, our experience has shown that the results are however hard to interpret, even by an expert.

280 For this reason, for the sake of interpretability, in our methodology we proceed by restricting the set of features for
281 further analysis. Hence, in the next step of the methodology, we select a small subset $R \subset [1, k]$ of features. The set R
282 should contain only non-redundant features, and should contain those that are highly informative with respect to the
283 set \mathbf{y}^A . The process we apply to select R is as follows. First, rank the features of \mathbf{x}^A by the amount of information they
284 provide about \mathbf{y}^A . Then, using this ranking, select the top features in order, removing those that are redundant with
285 respect to the previously selected features. This process stops when a certain number of features have been selected,
286 or when the next feature in the ranking gives little information about \mathbf{y}^A .

287 3.6. Clustering Using the Most Relevant Features

288 The next step of the methodology is to use the features in R to classify the anomalous scenarios into meaningful
289 types. We have observed empirically that the features in R naturally belong to different network aspects (e.g., radio
290 features, TCP features, QUIC features, etc.), and that an anomalous scenario may show issues in several of these
291 aspects simultaneously. This leads to a process in which subsets $r \subseteq R$ are selected, and each subset is used to classify
292 the anomalous scenarios using binary clustering indicating whether the selected subset of features seems to be (at
293 least partially) responsible for the anomaly. The conjecture is that this (unsupervised) process will select subsets that
294 correspond to different network aspects, and the clustering will separate “good” scenarios from those with issues in
295 that aspect. While this process is unsupervised, our experience is that it leads to subsets of R that an expert can easily
296 match with particular aspects of the network.

297 Hence, in this step we select α subsets $r_a \subseteq R$, $a \in [1, \alpha]$, and apply each of them to divide the set A of anomalous
298 scenarios into two clusters, a *low-performing* cluster (or a cluster with issues with respect to that network aspect,
299 indicated with a bit 0) and a *high-performing* cluster (bit 1). The bit labels 0 and 1 are assigned automatically to
300 the clusters, depending on which of them has better performance (in terms of QoS). Moreover, it is desirable to have
301 balanced, compact and non-redundant clusters. Each anomalous scenario will be assigned to one of the two clusters
302 generated with each subset r_a . This can be represented as a binary string of length α , with a value 0 or 1 in each
303 position, which classifies all anomalous scenarios A into 2^α types. Each of the bits in the vector associated to a
304 scenario conveys whether that scenario shows issues in the particular network aspect.

305 3.7. Aspect Classification of Scenarios

306 Finally, a second classifier C_2 is built using \mathbf{x}^A and \mathbf{y}^A as training data, and the 2^α types as the classes to which
307 these scenarios are assigned. Observe that we do not restrict the classifier C_2 to use only the features from R (although
308 it is very likely that they will be used). Thus, C_2 uses for training \mathbf{x}^A , \mathbf{y}^A , and the α -bit string of every scenario in A
309 (identifying the presence of anomalies in the α selected network aspects, as described in Section 3.6). The output of
310 C_2 will group anomalies into 2^α classes, each with a specific anomaly pattern.

311 It is important that this classifier is interpretable, because this makes experts able to understand why an anomalous
312 scenario has to be considered to have issues or not in each of the automatically identified network aspects. Note
313 that if an expert had been involved in the most relevant feature selection process, she should have assigned a specific
314 network aspect to each subset r_a , $a \in [1, \alpha]$, and possibly a meaning to every binary value of every element of the
315 network aspect anomaly string. Be it the output of manual inspection or automatic data analysis, this information can
316 be used to determine the issues each particular anomalous scenario has, and can be sent to the appropriate department
317 to promptly intervene.

318 3.8. Using the Classifiers to Detect and Identify Anomalies

319 After the previously described steps have been completed, we found ourselves with two classifiers C_1 and C_2 ,
320 that can be used to detect anomalous scenarios, determine network aspects that make this scenario anomalous, and
321 hence notify about this to the appropriate department for troubleshooting. Now we can use these trained classifiers to
322 analyze fresh data inputs and detect anomalies. In particular, consider a new scenario $(x_1, \dots, x_k, y_1, \dots, y_\ell)$. Using

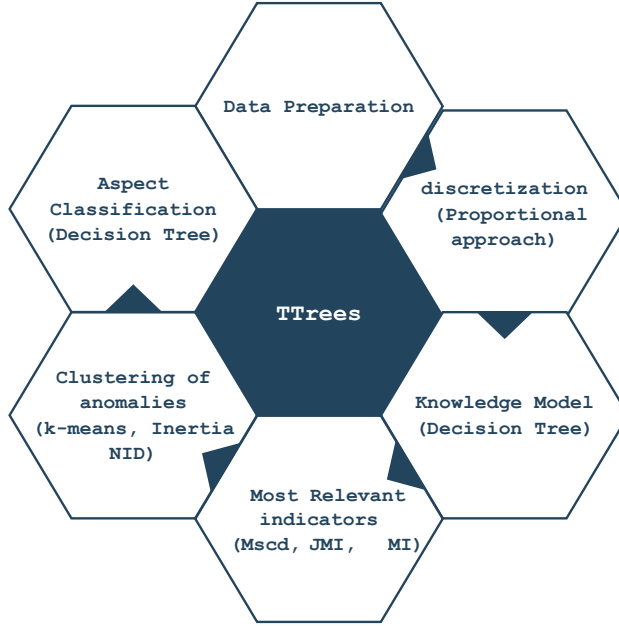


Figure 1: Steps of TTrees (using multiple ML techniques): beginning with data preparation, proportional discretization, training of knowledge tree, selection of the most relevant features, identification of anomaly clusters, and training of a network aspect anomaly classifier.

323 classifier C_1 we determine that this is an anomalous scenario if the KPI class C_1 assigned to (x_1, \dots, x_k) is not the one
 324 corresponding to the discretization class of (y_1, \dots, y_ℓ) but it is actually better. If that is the case, C_2 will be used, so
 325 that the class assigned to $(x_1, \dots, x_k, y_1, \dots, y_\ell)$ is an α -bit string expressing whether this scenario has issues in each
 326 of the α different network aspects or not, which allows for alerting the appropriate department(s).

327 4. Implementation of the CIAN methodology

328 Here we present the implementation details of the tool TTrees that automatizes CIAN, the methodology presented
 329 in the previous section. Fig. 1 depicts the different phases of the proposed methodology according to their implemen-
 330 tation in TTrees. As reported in the figure, our implementation leverages a number of standard tools for data analysis
 331 and ML, from discretization algorithms to clustering and classification.

332 4.1. Data Preparation in TTrees

333 The first phase consists in preparing the available data. This is a necessarily preliminary step that includes data
 334 filtering, cleaning of the dataset from the presence of incomplete entries, and extracting the k features in \mathbf{x} and the
 335 target KPIs in \mathbf{y} that need to be observed (i.e., note that the given dataset might include more features and KPIs than
 336 what we are interested in). Thus, we use a Python script to extract \mathbf{x} and \mathbf{y} from the dataset. The script automatically
 337 discards entries in which one or more KPIs or features are not reported. It also discards features that are constant
 338 through the dataset, which would not be informative. Features can be either numeric values or categorical entries
 339 (e.g., labels used to indicate the operational conditions of the network, like the name of the radio technology used, the
 340 name of the protocol adopted for transmission, etc.).

341 4.2. Discretization in TTrees

342 The selected target \mathbf{y} is usually a set of one or more numerical KPIs, whose values may lay on a continuous real
 343 interval. Thus, a first problem to address in the training of TTrees consists in discretizing the continuous target into
 344 a finite number of classes (i.e., categories). Using a simple quantization of the values in \mathbf{y} would introduce an error
 345 in the subsequent data analysis. Hence, instead we keep the target values as they are in the dataset, and assign a label

346 to each point, which identifies their class. Since TTrees is unsupervised, discretization uses progressive numbers as
347 labels (e.g., “0” to “5”).

348 In our methodology, being based on the study of the compressibility of the training dataset using optimal lengths
349 codes, it is of utmost importance that the discretization algorithm applied does not alter the distribution of the training
350 samples, that is, both the continuous KPI and the discretized version should follow the same distribution. A uniform
351 discretization is an easy to implement approach that satisfies this requirement. Therefore, we adopt a uniform dis-
352 cretization of target KPIs.

353 The number of classes (labels) used is not pre-defined. Instead, it is automatically derived by TTrees from the
354 available number of points in the dataset. There exist a collection of candidate discretization algorithms that are
355 not biased and have low variance (e.g., equal width, equal frequency, or fixed-frequency). Unfortunately, they require
356 optimizing hyperparameters [33], which makes them inadequate for an unsupervised automated software tool. Instead,
357 we use a discretization approach that does not require any hyperparameter tuning: *proportional discretization* [34].
358 This method uses a number of categories proportional to the size of the dataset. The number of categories used,
359 denoted by m , is computed as $m = (\log_2 n)/2$, where n is the number of samples of KPIs \mathbf{y} . We use the proportional
360 discretization approach to identify the centroids of the m intervals of KPI values (associated to the m categories). After
361 identifying the centroids, we apply k -means clustering to the continuous target variable \mathbf{y} , which assigns labels to the
362 data points and returns the boundaries between categories intervals.

363 4.3. Knowledge Model for Identifying Anomalies

364 Once the data points \mathbf{y} are assigned to a finite number of categories, TTrees continues by building a knowledge
365 model, meant to identify anomalies. For this step we use an ML classifier, which we want to be interpretable, so as to
366 allow the user to understand what kind of anomaly is detected. This is achieved by training a decision tree [35] with
367 the data set \mathbf{x} , and using as classes the categories defined in the discretization phase. The decision tree infers the class
368 for a sample \mathbf{y}_i of \mathbf{y} from the values of the associated features \mathbf{x}_i . This justifies the use of terms “knowledge model”
369 and “knowledge tree” for the role of this classifier. However, some samples \mathbf{y}_i can be misclassified, representing
370 anomalous behavioral trends that cannot be learned. Note that, with more advanced models (neural networks, random
371 forests, etc.) instead of a decision tree, we could gain accuracy, but we would lose interpretability and the model may
372 discard interesting samples.

373 The knowledge tree is build by applying the widely adopted Classification And Regression Tree (CART) algo-
374 rithm [36], with the Gini impurity metric [37] computed by comparing the output of the classification with respect to
375 the discretization categories. We refer the resulting decision tree as C_1 . The depth of C_1 is limited to a maximum of
376 $\lfloor (\log_2 n)/2 \rfloor$, so that, on average, the leaves of the tree hold a sufficiently large number of samples to properly classify
377 the KPIs. A deeper tree would suffer a high risk of overfitting [38]. Moreover, we also limit the number of samples at
378 the tree leaves to a minimum of five, to further reduce the risk of overfitting [37].

379 After the full tree has been derived to the maximum depth allowed, we perform a cost-complexity pruning, which
380 is the best known approach to avoid overfitting in potentially large trees [37], combined with cross-validation of the
381 candidate pruning points. The leaves of the tree that have the highest Gini impurity measure are removed first. Pruning
382 minimizes a cost-complexity metric that linearly combines the classification error (cost) made by the tree and the size
383 of the tree (complexity). However, the pruning metric has a hyperparameter which represents the relative importance
384 of cost with respect to complexity. Cross-validation is used to tune that hyperparameter automatically. In practice,
385 for each candidate value of the hyperparameter, we identify the tree with the minimal cost-complexity, and cross-validate
386 the performance of that tree by splitting the dataset into smaller pieces of at least 30 samples each: we use all but
387 one of the pieces to train the tree and the remaining one for evaluation. This process is repeated as many times as the
388 number of pieces in which we break the original dataset (so as to use every piece for validation exactly once).

389 Once C_1 is created, for each sample of \mathbf{y} we have two possibilities: the sample is classified or not under the same
390 class in the discretization phase and by the knowledge tree. If the knowledge tree returns a class higher than the one
391 of the discretization phase, an anomaly of interest is identified and included in set A .

392 4.4. Most Relevant TTrees Features

393 The next step consists in identifying a subset of features that are relevant for the anomalous scenarios in A , so as
394 to identify which network aspects are relevant to explain the anomaly. TTrees then evaluates each of the available

395 features by computing how much information on the target \mathbf{y} is contained in the feature. For this task we use the
 396 well known metrics *mutual information* (MI) [39] and *joint mutual information* (JMI) [40], from classic information
 397 theory.

398 We have also tested a novel metric that we have built on top of the concepts of Kolmogorov complexity and
 399 normalized compression distance (NCD) [41], which we call *miscoding* (Mscd). The miscoding of a sequence of p
 400 different features $\mathbf{x} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p\}$ with respect to the target \mathbf{y} is defined as

$$\text{Mscd}(\mathbf{x}_i, \mathbf{y}) = \frac{1 - \text{NCD}(\mathbf{x}_i, \mathbf{y})}{\sum_{j=1}^p (1 - \text{NCD}(\mathbf{x}_j, \mathbf{y}))}. \quad (2)$$

401 We propose Mscd because it measures the difficulty of reconstructing the target \mathbf{y} from the features of \mathbf{x} and vice versa,
 402 which accounts for the redundancy in \mathbf{x} . In TTrees we start by computing a conditional redundancy matrix with the
 403 normalized compression distance of all possible pairs of features with respect to the target variable $\text{NCD}(\langle \mathbf{x}_i, \mathbf{x}_j \rangle, \mathbf{y})$,
 404 based on a joint discretization of the attributes, and the computation of optimal length codes, given the relative fre-
 405 quencies of the discretized vectors. Then, we select the attribute with the highest NCD, and recompute the values
 406 of the redundancy matrix assuming that this value is selected. We repeat this process of selection of the best feature
 407 and recalculation of the redundancy matrix until all the features have been selected. The final miscoding is given by
 408 normalizing over one minus the NCDs computed for the different attributes.

409 MI is applied on a per-feature basis and is able to quantify the relevance of a feature, but it fails to identify
 410 redundancy. JMI is instead used on feature pairs, so that it allows to evaluate not only the relevance of an feature but
 411 also the redundancy with another feature. However, JMI is prone to errors in case outliers are present. Mscd offers
 412 the possibility of evaluating relevance and redundancy, plus the quantity of irrelevant information contained in the
 413 feature.

414 We will compare the performance obtained by applying MI, JMI and Mscd in the numerical evaluation section.
 415 Here it is enough to mention that in all of the three cases, we sort the features in decreasing order according to the
 416 selected metric, and pass the top of the list to the next TTrees phase. In particular, since we will need to cluster the
 417 anomalies based on selected network aspects, we pass to the clustering phase a number of features much larger than
 418 the number of aspects to study. For this paper we have used m^2 features, where m is the number of bins in which the
 419 target KPI has been discretized.

420 4.5. Clustering of Anomalies in TTrees

421 Each *pair* of relevant features is regarded as a potential *network aspect*. For each aspect, TTrees applies a clus-
 422 tering algorithm on all anomalous samples. This bidimensional clustering is done with the k -means algorithm [37],
 423 using the normalized values of the features obtained via RobustScaler [42]. In addition, and for visualization purposes
 424 only, when we have just one target KPI, we use a simple regression with respect to \mathbf{y} on each of the features, and we
 425 sort samples according to increasing regression values. This allows to plot anomalous samples consistently so that
 426 higher \mathbf{y} values are at the top-right corner of the plots. As a consequence, the cluster of samples at the top-right corner
 427 contains the samples with highest \mathbf{y} values and the cluster at the bottom-left corner the samples with lowest \mathbf{y} values.
 428 If the target KPI \mathbf{y} is better when larger (resp., smaller), then the bottom-left (resp., top-right) cluster is the one with
 429 issues in the network aspect defined by the two features. In general, the cluster with issues is assigned a label 0 in this
 430 aspect, and the other cluster is assigned 1.

431 The problem is that we have potentially a large number of pairs of features (and hence network aspects), and some
 432 of them might bring redundant information. Therefore, TTrees selects only a few pairs of features as network aspects,
 433 so as to easily and interpretably identify non-redundant descriptions of anomalies. In our current implementation we
 434 select $\alpha = 4$ feature pairs, i.e., the 4 most relevant aspects to evaluate to understand the anomaly. Thus, the number of
 435 features passed from the previous step is $p = \alpha^2 = 16$. We do not use all possible features for a matter of practicality.
 436 For example, if we had $p = 120$ features, which is a reasonable value for cellular data traffic traces, we would have to
 437 evaluate $p(p - 1) = 14280$ pairs, and make hundreds of millions of comparisons to test the redundancy among pairs.

438 The most relevant and descriptive feature pairs are selected based on multiple metrics. First of all, TTrees com-
 439 putes the inertia score [43] of the clusters resulting from each feature pair. The inertia of a cluster quantifies the
 440 tightness of the clusters (the lower, the better). The resulting ranked list is filtered by using two criteria: TTrees only

441 keeps clusters with low redundancy and which are balanced in the number of elements in the two clusters. Redun-
442 dancy of clusters is computed using the normalized information distance metric (NID) [44]. Here the order matters,
443 since if two feature pairs yield redundant clusters, the feature pair with lower inertia is retained, while the other is
444 discarded. In order to automatically tune and apply both filters, we select filtering thresholds with the help of the
445 histograms obtained for the distributions of redundancy scores and balance metric. We have observed that the his-
446 togram of redundancy shows a bimodal distribution, therefore we select as threshold to accept/reject a feature pair the
447 point with the lowest value in-between the two peaks of the bimodal distribution. For balancedness, the distribution
448 histogram is multimodal, with an automatic threshold selection where the minimum point between first and second
449 peaks (minimum accepted value of balancedness), and the minimum point between the last two peaks (maximum
450 accepted value for balancedness).

451 4.6. Aspect Classification in TTrees

452 The α top feature pairs, obtained in the previous phase of TTrees, identify network aspects in which the anomaly
453 could be rooted. In this phase they are used to train an interpretable classifier (again a decision tree) to map the
454 anomalies A onto the 2^α classes corresponding to the combinations of the α network aspects relevant for the anomalies.
455 This *aspect classifier*, denoted as C_2 , uses the binary values of the aspects assigned in the previous phase, and builds
456 new categorical target variables consisting of binary strings that combine these binary aspect values. The elements
457 with value 1 in the binary string identify aspects in which the scenarios of that category may have issues, and for
458 which a network specialist should be called. Note that this classifier makes decisions based on the full list of features,
459 and not only based on the most relevant features identified by TTrees. In the training of C_2 , like for C_1 , we avoid
460 overfitting by limiting the depth of the tree to $\lfloor (\log_2 |A|) / 2 \rfloor$, and apply cost/complexity pruning with cross validation
461 to increase the generalization level of the tree.

462 4.7. Software Implementation

463 TTrees is implemented using Python. We use two libraries: (i) the widely adopted scikit-learn library, which
464 provides us with ML algorithms, such as decision trees and k -means, and (ii) a library developed by us for automatic
465 ML tools, which includes an Mscd calculator.

466 5. Empirical Evaluation

467 Existing troubleshooting tools do not generalize very well and need the assistance of an expert. For this reason,
468 here we evaluate TTrees along with a supervised version of our tool, which provides a validation baseline. For
469 simplicity, here we use only one target KPI (e.g., the throughput) and restrict the analysis of anomalies to the case
470 in which the knowledge model predicts better throughput than what observed. We also present a modified version
471 of TTrees, in which we modify the clustering algorithm in order to effectively analyze imbalanced datasets. The
472 first datasets used in this section were extracted from experiments executed in controlled environments. These are
473 complemented with Google and Facebook datasets, which are the result of homogeneous cellular network conditions
474 and include nodes deployed in both static and mobile stations. The latter can suffer from network outages and undergo
475 unexpected variances in the quality of the signal, which enriches the possible scenarios captured in the dataset and
476 can be a clear source of anomalies. As it will be later developed, these uncertain conditions yield aspect classification
477 clusters where the number of anomalies is imbalanced, that is, the ratio is not equal between clusters and one class
478 accumulates most of the samples.

479 The datasets collected for this paper and the analysis of all datasets presented here, for all services and operators
480 discussed, are available in GitHub.⁵

481 5.1. Supervised ML-based Troubleshooting

482 The supervised version of TTrees, denoted as STrees, uses a supervised selection of network aspects that are most
483 relevant to identify and cause anomalies in the problem classification phase of the methodology. In particular, with
484 the supervision of a network expert, STrees selects as network aspects those combinations of features that are most
485 relevant for the selected target KPIs, as identified by the expert. Afterward, in STrees, we take the tagged classes and
486 use them to train C_2 , i.e., the final aspect classification tree, like is done in TTrees.

⁵<https://github.com/Mohmoulay/WoWMoM2021>

Table 2: Brief description of the experimental datasets used in this work

ID	Dataset	# Samples	# features	Families of features
#1	HTTP FILE DL	6,690	228	Radio, RTT, Duration, TCP Volume, TCP Flags, TCP Window, Packet Anomaly
#2	HTTP LIVEPAGE DL	10,500	126	Radio, RTT, Duration, TCP Volume, TCP Flags, TCP Window, Packet Anomaly
#3	HTTP FILE DL MONROE	120,000	106	Radio, RTT, TCP Window, Packet Anomaly
#4	QUIC FILE DL MONROE	3,951	91	RTT, Duration, TCP Volume, TCP Flags, TCP Window, Packet Anomaly
#5	Facebook Page Download	1447,839	109	Radio, RTT, TCP Window, TCP, Packet Anomaly
#6	Google Page Download	1261,119	109	Radio, RTT, TCP Window, TCP, Packet Anomaly

5.2. Datasets

Table 2 summarizes the characteristics of the datasets used in this work. We first analyze two data collections of cellular network experiments used by Nokia in 2019 for network performance assessment of various European 4G networks (Datasets #1 and #2). These two datasets contain drive-test information on TCP traffic experiments in which constant-size (3 MB) files are downloaded, and live web-pages are fetched with a mobile user device. Then, we further evaluate TTrees using larger mobile broadband connectivity datasets collected with MONROE [45]. This includes generating a dataset of QUIC *qlogs* by running experiments from MONROE nodes.

More in detail, the Nokia datasets used for experimental validation are real test records generated by processing all the information provided by the testing equipment used in the drive-tests, and aggregating the information at test level (count, sum, min, max, average, percentiles, etc.). The resulting datasets have a single row per test and hundreds of columns summarizing all the dimensions (i.e., context features like date, time, location, network element information, etc.), technology features (radio, TCP/IP, application, etc.) and the KPI related to that particular test. The data transmitted in the drive tests is synthetic and does not include customer’s sensitive information, protected by the European Union General Data Protection Regulation (GDPR) or similar regulations (i.e., the data has been generated by a testing device and not by real users). Finally, potentially sensitive data, such as the identity of the network operator, has been anonymized. The Nokia datasets utilized in this study are rich but not sufficiently large to, e.g., train a neural network. They contain 54 228 and 21 655 rows, respectively, with 1326 columns.

The third dataset we use (Dataset #3) is a MONROE dataset with 120 000 rows and 106 columns. We have obtained access to MONROE and replicated the experiments of the first Nokia dataset (i.e., with 3 MB HTTP file downloads) on a larger scale, although with a different number of features. Moreover, some of the MONROE dataset highlights include the ability to filter by test type, infrastructure, operator, vendor, and transmission technology.

The fourth dataset (Dataset #4) was also obtained by means of the MONROE platform with the main aim of capturing traffic different than TCP. The experiments conducted in this scenario consist in the exchange of data packets from a server we control to a client residing on a MONROE network probe. Packets are sent over QUIC. From a data transport point of view, the novelties introduced by QUIC can be roughly summarized as follows: QUIC offers a combination of TCP reliability and multiplexing capabilities with the speed provided by the UDP transport protocol, over which QUIC is built [46]. In the experiments, we tested multiple possible variants for the congestion control algorithm implemented by QUIC, those being BBR, COPA, New Reno and Cubic. We also conducted two classes of QUIC experiments with sequential and parallel file downloads, respectively. The server can utilize multiple streams to exchange data with the client over QUIC. The number of streams spawned ranged from 1 to 5 across the experiments conducted. To limit the duration of the experiments, each opened connection had a lifespan of 15 seconds, enough to evaluate the behaviour of congestion control algorithms over variable radio conditions. Client nodes located in Norway, Sweden and Spain were utilized to run the experiments. Two types of clients were identified according to their mobility: static and mobile nodes. Among the mobile nodes available in the MONROE platform, the experiments conducted involved those which were installed on Swedish buses. Introducing mobile nodes contributes to increase the heterogeneity of samples with scenarios where the quality of network connectivity can be significantly lower than

523 in static cases.

524 The complexity of the QUIC protocol and its use of encryption can hinder the inspection of QUIC traffic to such
525 an extent that a specific format called *qlog* has been defined to organize and collect the information captured during
526 a data exchange [47]. A *qlog* file includes not only QUIC events, but also protocol events, each of which includes a
527 timestamp, the type of event and the value associated as specified in [48]: the timestamp indicates the time at which
528 the event was registered, a field describes the type of event associated to the record, and finally, the data recorded
529 in included. These *qlogs* have been the primary source of information to produce the data that populates the fourth
530 dataset, which emulates the format and features present in the Nokia dataset. Therefore by applying TTrees to QUIC
531 *qlogs*, we test its flexibility to work with heterogeneous kinds of data without the need of tuning any hyperparameter
532 manually.

533 Datasets #5 and #6 were specifically designed to evaluate the connectivity to websites that should offer high
534 performance, so that possible network anomalies could be experienced mostly because of radio problems. In this way,
535 we obtain datasets in which radio and transport anomalies are strongly imbalanced. Specifically, Dataset #4 collects
536 measurements run against the Facebook webpage. We obtained the measurements by means of multiple Facebook
537 webpage downloads in various European countries, using MONROE probes. The average downloaded webpage size
538 of Facebook was 72 757.3 bytes and, at each experiment run, we cleaned the cache to make sure that a new complete
539 download occurs. In some countries, such as Sweden and Norway, we had access to mobile nodes mounted on public
540 transport vehicles. Thus, the downloaded page statistics for those countries are a mixture of statistics for static and
541 mobile cellular nodes, hence the slightly higher radio problems which is depicted later in this subsection. As for Italy
542 and Spain, we only report results for static cellular nodes. We run similar measurements with MONROE and Tstat
543 against the Google search engine webpage, which was used to build Dataset #6. In those experiments, We observed
544 an average downloaded size of 17 368.8 bytes, following the same procedure as the Facebook service (repeating
545 downloads, clearing the cache, etc.).

546 In what follows, we use the first Nokia dataset to validate CIAN/TTrees and showcase its troubleshooting capabil-
547 ities for the case of HTTP file download operations over LTE networks. We only consider data for a single operator
548 (we anonymize the operator’s name for privacy purposes). The other datasets are used to showcase the potentials of
549 CIAN and pros and cons of how the methodology is implemented. In particular, the second Nokia dataset is for HTTP
550 webpage download performed live, over LTE networks, and one single operator as well. The first MONROE dataset
551 Dataset #3 is for HTTP file download over LTE networks, and with one single operator, with tests anonymously
552 conducted from labs and public buses in various European countries; so it is similar to the first Nokia dataset, but
553 with more entries, less features and much more geographical diversity. The MONROE QUIC dataset Dataset #4 is
554 relatively small, with less than 4000 entries, and reports very different families of features with respect to the other
555 datasets. Datasets #5 and #6 are used to show that, in some cases, we need to trade off accuracy for interpretability of
556 CIAN results.

557 5.3. TTrees in Action and its Validation

558 For the analysis of the HTTP file download dataset of Nokia, we choose the *throughput* as the desired target KPI.
559 Since there are 6690 samples, the number of categories used for this specific case is $\lfloor (\log_2 6690)/2 \rfloor = 6$. Fig. 2
560 illustrates the results obtained when applying the proportional discretization approach to this KPI in TTrees. The
561 discretized target variable is fed to the CART algorithm to build the knowledge tree. In CART, we fit the knowledge
562 tree using all the 228 features available in the dataset.

563 Fig. 3 explains a subset of the resulting knowledge tree. For a given sample at the root of the node, the CART
564 algorithm checks if the *Abs.RWIN.FirstSec* feature (the absolute value of the *TCP received window* during the first
565 second of a connection) is less or equal to the self-learned threshold of about 1.3 MB. If this statement is true, the
566 CART algorithm moves to the left side of the tree. Each node of the knowledge tree reports the value of Gini impurity,
567 i.e., a measure of the fraction of misclassified data points that are reported in the tree leaves that can be reached
568 from that node. It also reports the total number of samples examined in the node (6690 at the root) and the number
569 of samples that corresponds to each discretization category. E.g., [2185, 1646, 1293, 950, 506, 110] in the root node
570 correspond to the number of samples that fall in the “Very bad” (0), “Bad” (1), “OK” (2), “Good” (3), “Very Good”
571 (4), and “Excellent” (5) discretized throughput categories. The class value associated to each node is the one with
572 the highest number of samples (“Very bad”, in the example). The full knowledge tree generates a meaningful set
573 of classification rules, each based on a feature, and whose family (radio, TCP window, TCP flags, etc.) is reported

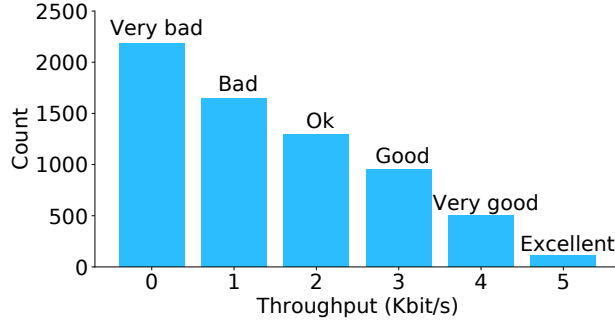


Figure 2: Target KPI (throughput) distribution for Dataset #1 via the proportional discretization approach

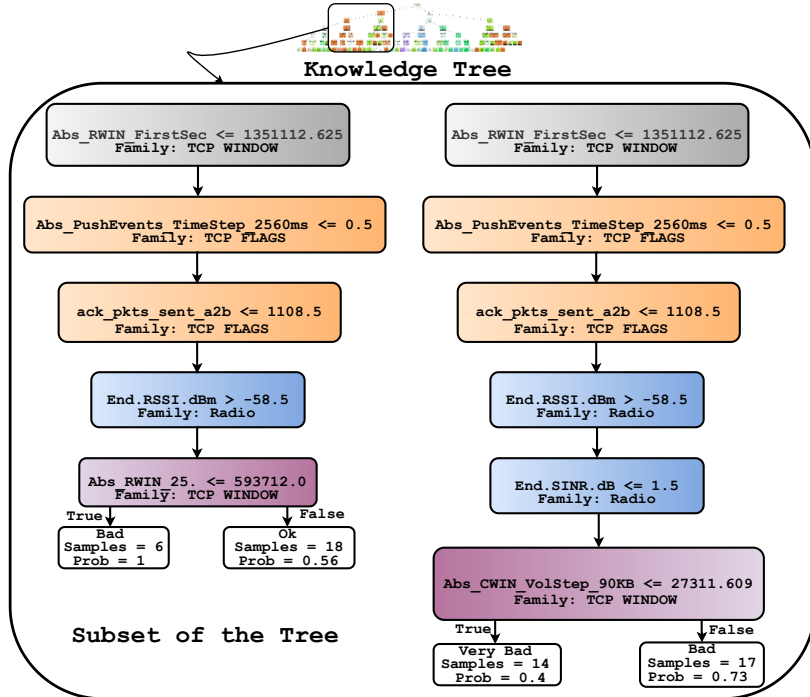


Figure 3: Zoom into a subset of the knowledge tree built for Dataset #1. The figure also reports which family of features the branching variable belongs to (see Table 2). Observe that both branches are the same up to the fourth node, after which the left corresponds to the *true* branch and the right to the *false* branch.

574 in each node of the tree. That can help expert troubleshooters get an initial idea of what is going on in the cellular
 575 network, but since the full knowledge tree is wide, manual inspection and interpretation can be cumbersome.⁶

576 Fig. 4 shows that the knowledge tree classifies the target samples with high accuracy, although the number of cases
 577 in which the predicted class is better than the one observed in the discretization phase is not negligible ('-1' in the
 578 figure). TTrees selects these 740 anomalies to continue the analysis.

579 The number of network aspects that can be analyzed by composing binary clusters in the aspect classifier C_2 with
 580 a population of 740 anomalies is $\alpha = \lfloor (\log_2 740)/2 \rfloor = 4$. Thus, the number of features to use is $\alpha^2 = 16$. Fig. 5 depicts
 581 the 16 most relevant features obtained by using the Mscd metric (left subplot) and by the MI and JMI metrics (right
 582 subplot). The figure clearly shows that MI and JMI practically find the same set of features, while Mscd identifies
 583 a significantly different subset, which is due to the fact that Mscd accounts for redundancy of features while MI and
 584 JMI do not.

585 In Figs. 6 and 7(a) the redundancy of the binary clusters obtained with the 16 most relevant features is shown,

⁶The full tree is large and cannot be correctly visualized here, so we made it available online at <https://github.com/Mohmoulay/WoWMoM2021>.

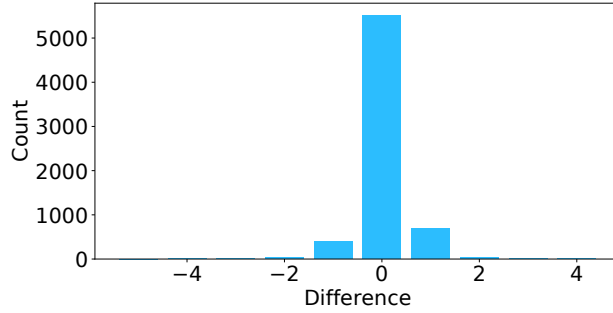


Figure 4: Difference (predicted - discretized) between the class assigned by the knowledge tree and the discretized category for the target KPI of Dataset #1

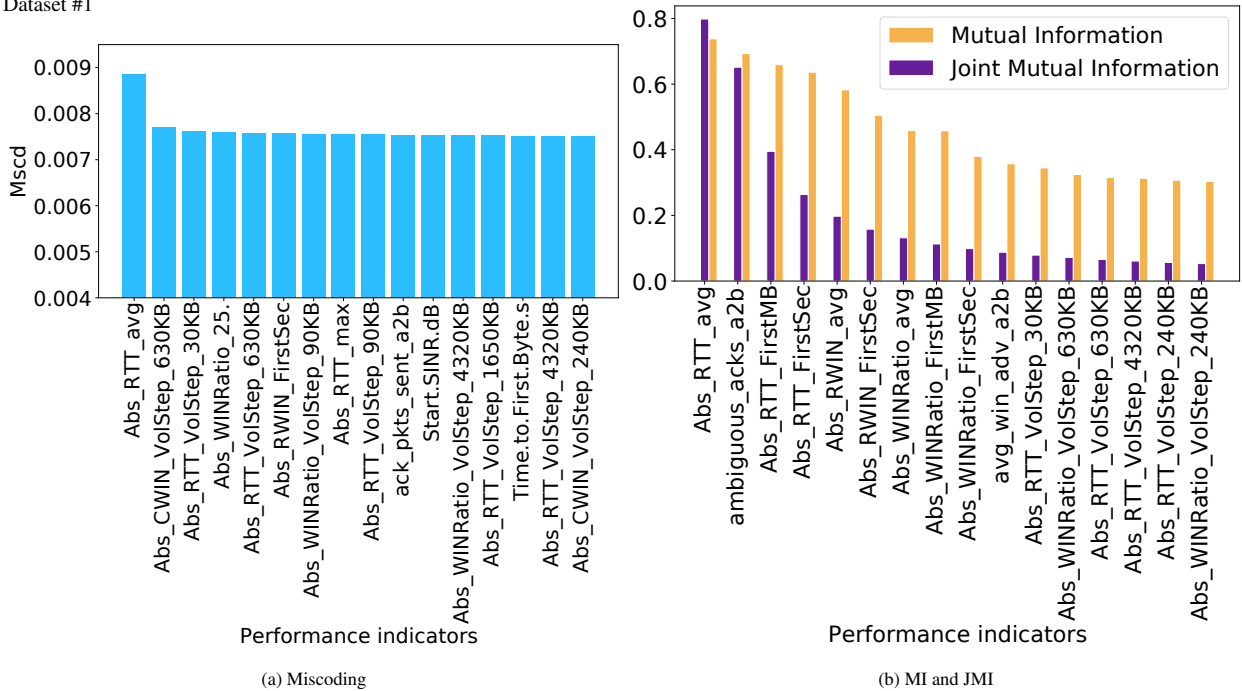


Figure 5: Ranking of features according to their relevance to the description of anomalies for Dataset #1

586 for MI, JMI, and Mscd metrics, respectively. The heatmaps shown in the figures represent the NID metric between
 587 cluster pairs: the darker the color the lower the distance between cluster pairs (and hence the higher their redundancy).
 588 Fig. 7(a) shows lighter colors than the other two heatmaps of Fig. 6, which means that Mscd yields features whose
 589 combinations of clusters have lower redundancy. For this reason we will mainly consider the features selected with
 590 Mscd in the rest of this section.

591 Fig. 7(b) illustrates the inertia of the clusters obtained with the features selected with Mscd, and the filtering
 592 process based on balancedness and low redundancy criteria. The output of this filtering process is a set of 4 network
 593 aspects (in practice, 4 pairs of features) used for binary clustering. The visualization of the filtering results for the
 594 selection of relevant aspects obtained with MI and JMI is omitted because it is not relevant for comparison purposes.
 595 What matters is instead the set of results on the performance eventually obtained with those metrics, which will be
 596 commented later and which are reported in Table 3.

597 Fig. 8 shows the 4 network aspects obtained with the first Nokia dataset and the corresponding clusters to be
 598 used in the aspect classifier C_2 . The figure shows that TTrees has identified anomaly clusters characterized by TCP
 599 operational parameters, radio parameters, delay parameters and packet anomaly. In STrees, instead, an expert would
 600 have to look manually at the “usual suspects” for the presence of anomalies, such as possible radio strength problems,
 601 TCP-related problem, packet anomaly and high RTT caused by buffering, and cluster them to see if they yield any

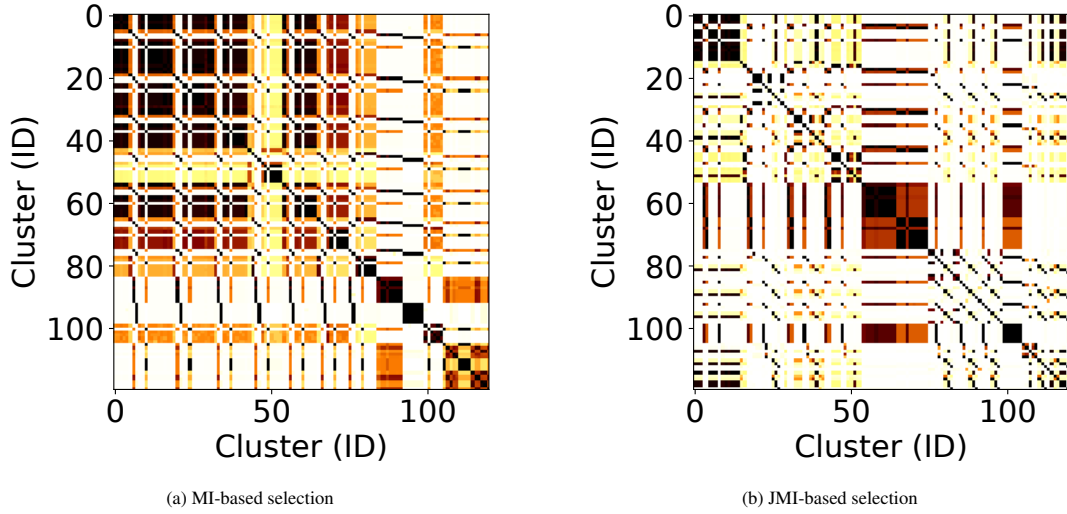


Figure 6: NID heatmap matrix for cluster pairs. Darker colors represent redundancy.

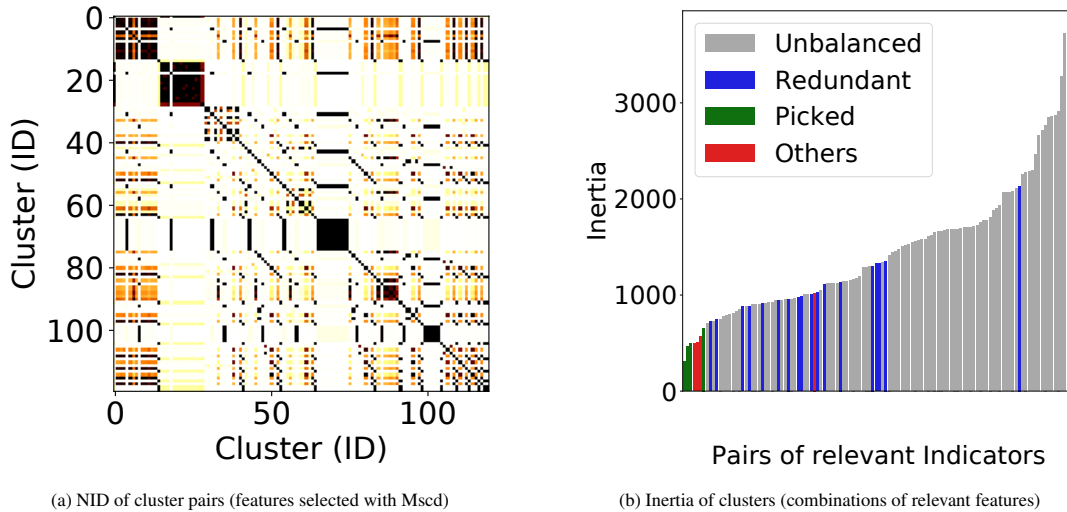


Figure 7: Network aspect selection based on miscoding; lower inertia of k -means clusters is preferred, unless clusters are redundant or imbalanced

602 groups. In this case, a Nokia cellular expert has identified for us the most relevant features, which led STrees to
 603 identify the network aspects visualized in Fig. 9. Aspects identified by TTrees and STrees are quite different, except
 604 in both cases they cover the same families of features (radio, TCP window, RTT, packet anomaly). The actual features
 605 selected by TTrees could be less intuitive for an expert (e.g., selecting number of acknowledgment instead of packet
 606 loss), but without loss in information. In reality, for an expert, Aspect 3 from Fig. 8 makes sense but may not be the
 607 first (intuitive) choice; instead, he/she would first look at the packet loss from Fig. 9. Indeed, one of the advantages
 608 of the TTrees methodology over STrees is the capability of creating meaningful yet not intuitive inter-family aspects
 609 combining radio with TCP, TCP window with packet anomaly, etc.

610 The corresponding C_2 trees trained with TTrees and STrees are (partially) reported in Fig. 10. The same four
 611 families of identifiers are deemed as relevant by supervised and unsupervised approaches: Radio, TCP-based, RTT,
 612 Packet anomaly. TTrees detects a dependency in Aspect 4 that is negatively correlated: when the value of the specific
 613 feature Abs_RTT_VolStep_630KB (i.e., the RTT observed after averaging over chunks of 630 kB of transmitted data)
 614 is higher than 87.909 ms, the throughput is going to be lower. This means that channel capacity will be the limiting
 615 factor causing buffering and we have a negatively correlated relation in Aspect 4. Now, if we look at the right side of

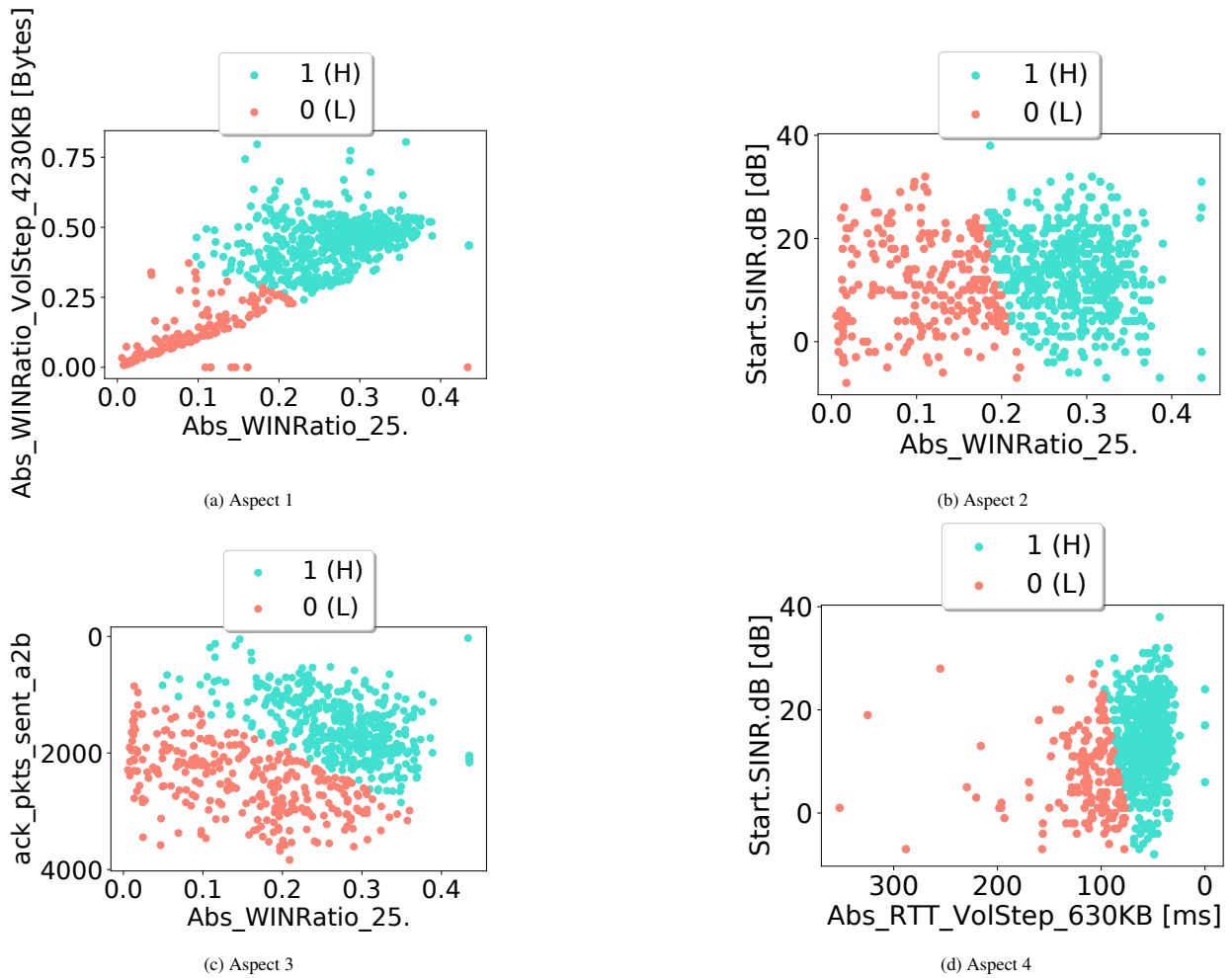


Figure 8: Aspects obtained with TTrees for Dataset #1 (using Mscd, see Figs. 5 and 7). The parameters used are Signal-to-interference-plus-noise ratio (SINR), Number of packet acknowledgment sent, absolute TCP Window ratio, and RTT.

616 Fig. 10, STrees is describing low RTT and high packet loss as well, deriving that result from the same TCP features
 617 plus radio. Furthermore, if we look at the whole tree, they both convey similar results. In practice, TTrees and STrees
 618 are comparable and show similar classification, taking into account the difference seen in the aspect selected. For
 619 instance, the root node is the same: Abs_RTT_avg (i.e., the average RTT observed during file download) is less or
 620 equal to 122.432 ms, but the tree leafs can not be identical. Both results are usable by an expert to identify manageable
 621 sets of similar tests for troubleshooting. Additionally, the tree rule would allow correct assessments to assign each
 622 group to the right troubleshooting department.

623 The overall accuracy score of C_2 for both approaches is 90%, which means that C_2 can identify 90% of non-
 624 anomalies and anomalies alongside their root cause. With the above simple example, we have defined an automated
 625 anomaly detection system that is functional to improve the quality of the networking service. Indeed, TTrees can
 626 self-identify network performance issues. So it can be used to alert the right support personnel, either the TCP or
 627 radio optimization team, in the presented example, which will receive as input the rules from Fig. 10 that raised the
 628 alarm. Acknowledging that other classification approaches might yield higher accuracy, we remark that would occur
 629 at the expense of interpretability.

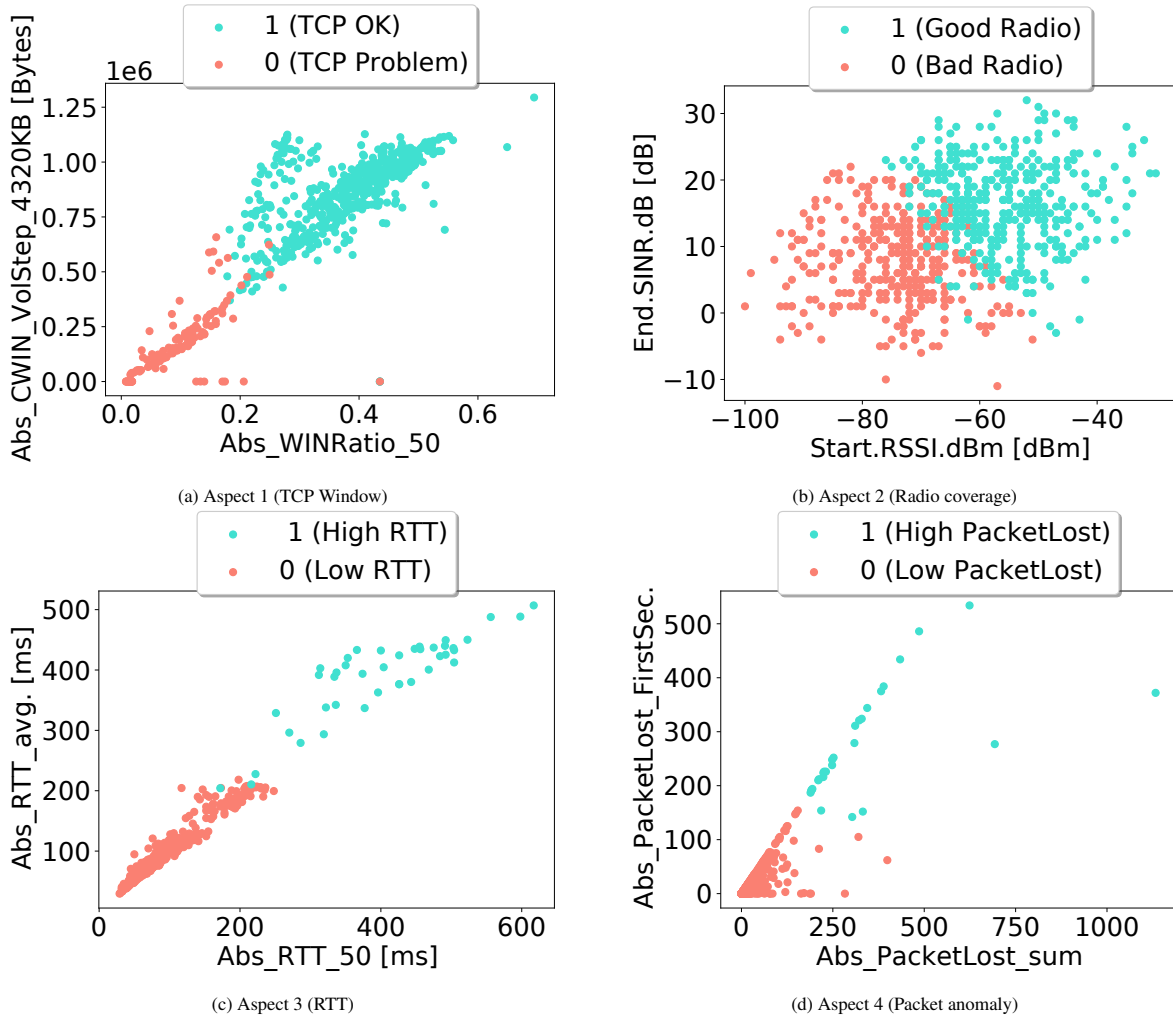


Figure 9: Aspects obtained with STree for Dataset #1

630 5.4. TCP Performance Evaluation

631 TTrees has been executed in three different datasets to identify anomalies with TCP data exchanges. The target
 632 KPI for Dataset #2 is *Average Session Duration*, discretized into six classes, and for Dataset #3 the KPI is *Throughput*,
 633 discretized into eight classes. Since the new datasets are large enough for TTrees, we split them into training and test
 634 subsets, with 10% and 20% testing data while applying cross validation 30 times, respectively. Note that, due to the
 635 small size of Dataset #4, data could not be split into training and testing sets and thus, cross validation could not be
 636 carried out, hence the “NA” entries on the table. Tests are done in classifiers C_1 and C_2 to get the metrics shown in
 637 Table 3. The number of anomalous scenarios is 1152 and 6902, with 25 and 36 most relevant features, respectively.
 638 TTrees with Mscd was able to identify similar network aspects in Datasets #2 and #3 as it was in Dataset #1: radio,
 639 TCP, and RTT. Unfortunately, the selection of the most relevant features did not work well using MI and JMI for
 640 Dataset #2 (HTTP livepage DL), since all the pairs of features selected were redundant, and TTrees could not build
 641 C_2 (hence the entries “-” in Table 3). This was not surprising to our expert, since the anomalous scenarios of Dataset
 642 #2 are very hard to classify, since they contain downloads of webpages with multiple sizes. The fact that Mscd found
 643 meaningful network aspects is indeed impressive. For Dataset #1 and #3, the C_2 classifiers show similar classification
 644 performance (see Table 3). However, the use of JMI and MI did not produce network aspects easily identifiable by
 645 an expert. In summary, using Mscd leads to better performance than the other two metrics, reaching good accuracy
 646 scores while producing meaningful sets of rules for troubleshooting.

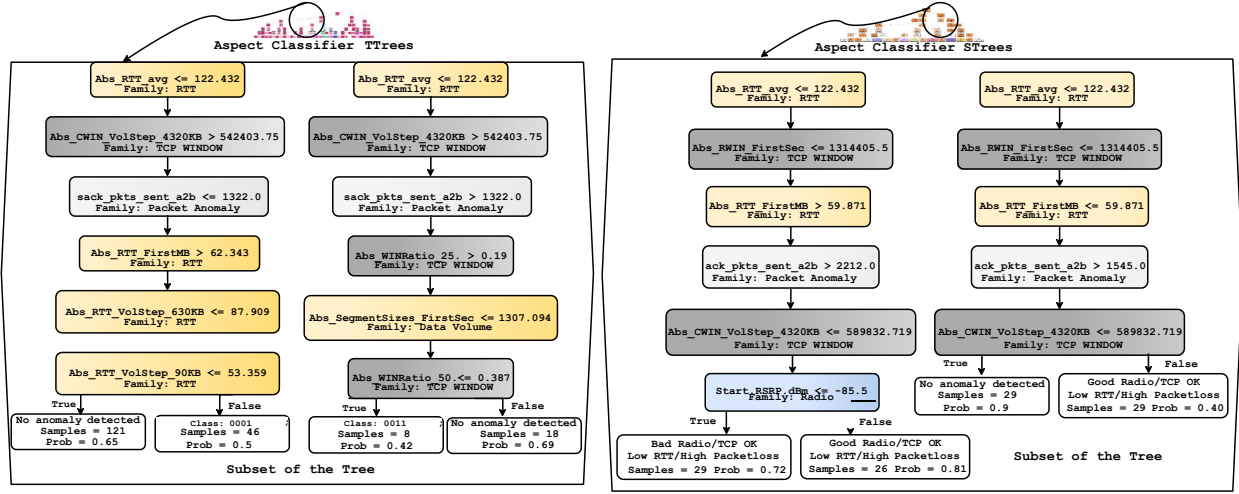


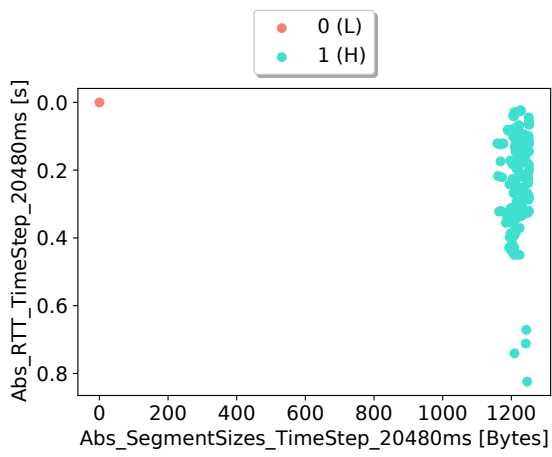
Figure 10: A comparison between the supervised and unsupervised approach showing the resulting aspect classifiers

Table 3: Knowledge (C_1) and Aspect Classification (C_2) trees' performance across the different datasets using multiple aspect selection algorithms (Mutual Information, Joint Mutual Information, and Miscoding)

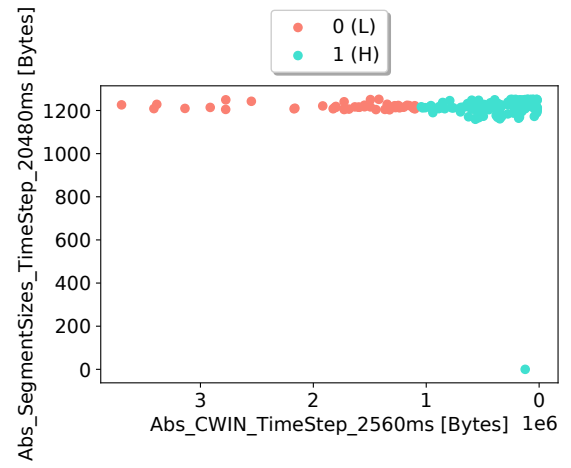
	Classifier	Accuracy	Precision	Recall	F1
Dataset #1 HTTP file DL (Nokia)	C_1	82.5%	NA	NA	NA
	C_2 (MI)	92.0%	NA	NA	NA
	C_2 (JMI)	91.0%	NA	NA	NA
	C_2 (Mscd)	90.0%	NA	NA	NA
Dataset #2 HTTP livepage DL (Nokia)	C_1	82.6%	80.0%	81.5%	80.0
	C_2 (MI)	-	-	-	-
	C_2 (JMI)	-	-	-	-
	C_2 (Mscd)	88.0%	85.0%	86.0%	84.0
Dataset #3 HTTP file DL (MONROE)	C_1	92.0%	90.0%	91.0%	89.0
	C_2 (MI)	95.0%	92.6%	94.4%	92.5
	C_2 (JMI)	95%	94.5%	94.5%	92.7
	C_2 (Mscd)	95.0%	93.0%	94.4%	93.5
Dataset #4 QUIC file DL (MONROE)	C_1	87.3%	NA	NA	NA
	C_2 (MI)	94.0%	NA	NA	NA
	C_2 (JMI)	94.0%	NA	NA	NA
	C_2 (Mscd)	96.0%	NA	NA	NA

5.5. QUIC performance evaluation

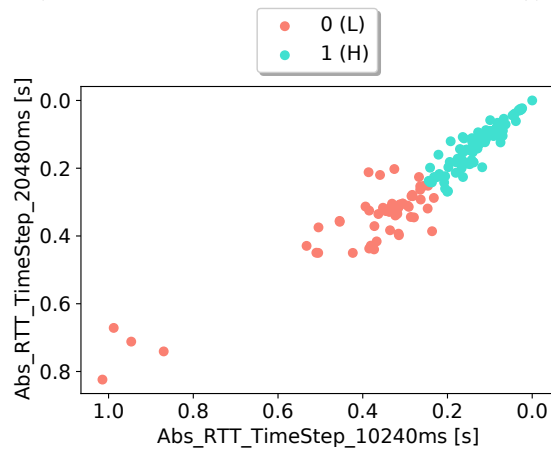
Finally, we use TTrees with our fourth dataset, containing the qlogs of almost 4000 experiments carried out with MONROE. The target KPI for Dataset #4 is the average throughput, which has been discretized into five classes. The congestion control algorithm and execution type used in each experiment were recorded jointly with qlogs and MONROE probe context parameters. Table 4 summarizes per-congestion control algorithm statistics computed to form the knowledge tree C_1 with the initial steps of the TTrees methodology. The table shows that experiments with COPA, Cubic and New Reno systematically yielded low throughput, especially with COPA, while with BBR we can observe a higher throughput, in terms of distribution over the 5 identified classes. The table also shows the number of anomalous cases classified according to network aspects. There were in total eight possible combinations of three network aspects in which to look for anomalies. The three identified network aspects are shown in Fig.11. The



(a) Aspect 1 (Segment sizes)



(b) Aspect 2 (TCP Window)



(c) Aspect 3 (RTT)

Figure 11: Aspects obtained with TTrees for Dataset #4

Table 4: Classification output of TTrees with QUIC experiments – Data grouped per congestion control algorithm (classes with no samples are omitted)

	Knowledge tree C_1		Aspect classification tree C_2	
	Throughput class	Count	Anomaly class	Count
BBR	0	363	111	926
	1	304	110	1
	2	243	011	29
	3	87	010	42
	4	1		
COPA	0	983	111	991
	1	14	011	6
Cubic	0	579	111	907
	1	170	011	9
	2	167	010	35
	3	42	000	7
New Reno	0	562	111	928
	1	212	011	16
	2	164	010	15
	3	40	001	3
			000	36

Table 5: Classification output of TTrees with QUIC experiments – Data grouped per experiment execution type (classes with no samples are omitted)

	Knowledge tree C_1		Aspect classification tree C_2	
	Throughput class	Count	Anomaly class	Count
Parallel	0	1260	111	1831
	1	413	011	40
	2	208	010	48
	3	70	001	3
			000	29
Sequential	0	1227	111	1921
	1	287	110	1
	2	366	011	20
	3	119	010	44
	4	1	000	14

657 first two of them show a clear dependency on one parameter only, which is the segment size in the first case and
658 the congestion window in the second case. Note that small segment sizes would naturally lead to low throughput
659 while small congestion window values are typical of bad performers. So, the selection of those network aspects,
660 which was done automatically by TTrees, makes sense. The third network aspect is a combination of RTT statistics
661 taken with different averaging intervals. Indeed, when it comes to evaluate the throughput of congestion-controlled
662 transmissions, it is well known that RTT statistics are a meaningful network aspect to consider as well. In particular,
663 high RTT values tend to indicate the presence of radio limitations.

664 Table 4 shows that not all combinations of network aspects were observed among the identified anomalies. Of
665 course, the more frequently observed case is that there are no anomalies (class 111), while most of detected anomalies

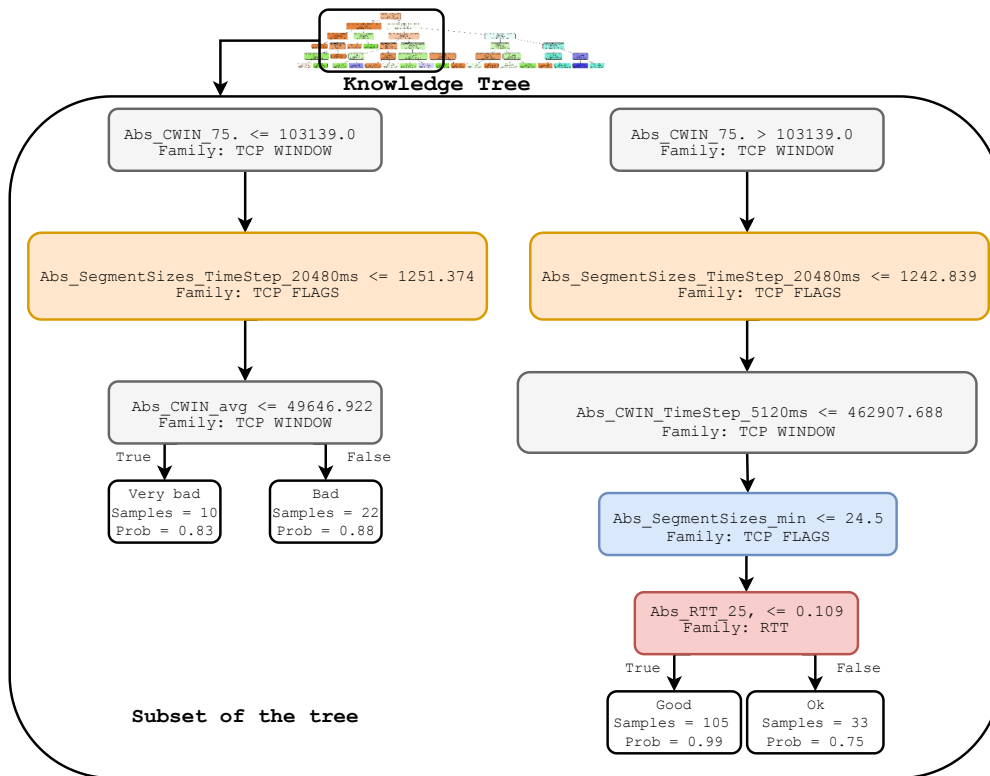


Figure 12: Zoom into a subset of the branches of the knowledge tree C_1 built by TTrees for Dataset #4. The figure also reports which family of features the branching variable belongs to (see Table 2)

666 show a combination of two or three network aspects. QUIC anomalies with BBR mostly depend on the segment
 667 size and RTT, while Cubic and NewReno anomalies also depend on the congestion window. The least number of
 668 anomalies was observed with COPA, and are caused by low values of the segment size utilized for transmission. Note
 669 that looking at Table 4, we can see that no experiment conducted using BBR falls under the second anomaly class,
 670 i.e., the congestion window network aspect does not affect BBR anomalies. This is a meaningful result since with
 671 BBR the value of the congestion window (CWIN) is automatically modulated to maintain the RTT within acceptable
 672 values, so that CWIN cannot be directly correlated with the experienced throughput.

673 If we consider parallel and serial QUIC data transmissions separately, Table 5 resumes the distribution of through-
 674 put samples—where of course we can see that sequential downloads observe higher per-download throughput—and
 675 anomalies. It is interesting to note that anomalies for parallel download experiments show a fundamentally different
 676 distribution across anomaly classes with respect to the case of sequential downloads. For instance, while the first
 677 network aspect (the segment size) is identified as a common cause of anomaly in all cases, the second network aspect
 678 (i.e., the congestion window size) is more relevant for parallel downloads. This result is meaningful because segment
 679 size always affects transmission rate, while the congestion window is affected by the number of ongoing transmission
 680 flows.

681 Note that, due to the small size of Dataset #4, data could not be split into training and testing sets and thus, cross
 682 validation could not be carried out. The knowledge tree, shown in Fig. 12, yielded 199 anomalous scenarios in total.
 683 As mention before, here we selected only three network aspects, which were extracted from a set of nine most relevant
 684 features. The aspect selection of TTrees automatically deemed appropriate to stop the search of network aspects to
 685 three, because adding a fourth one would have been redundant. We remark that this was done automatically, and we
 686 did not have to manually change the number of network aspects to select. The accuracy of the aspect classification
 687 tree for the QUIC dataset of Fig. 13, was the highest of all C_2 trees, with a value of 96% and includes the three types
 688 of anomalies identified. Also for the case of QUIC, using miscoding rather than (join) mutual information for aspect

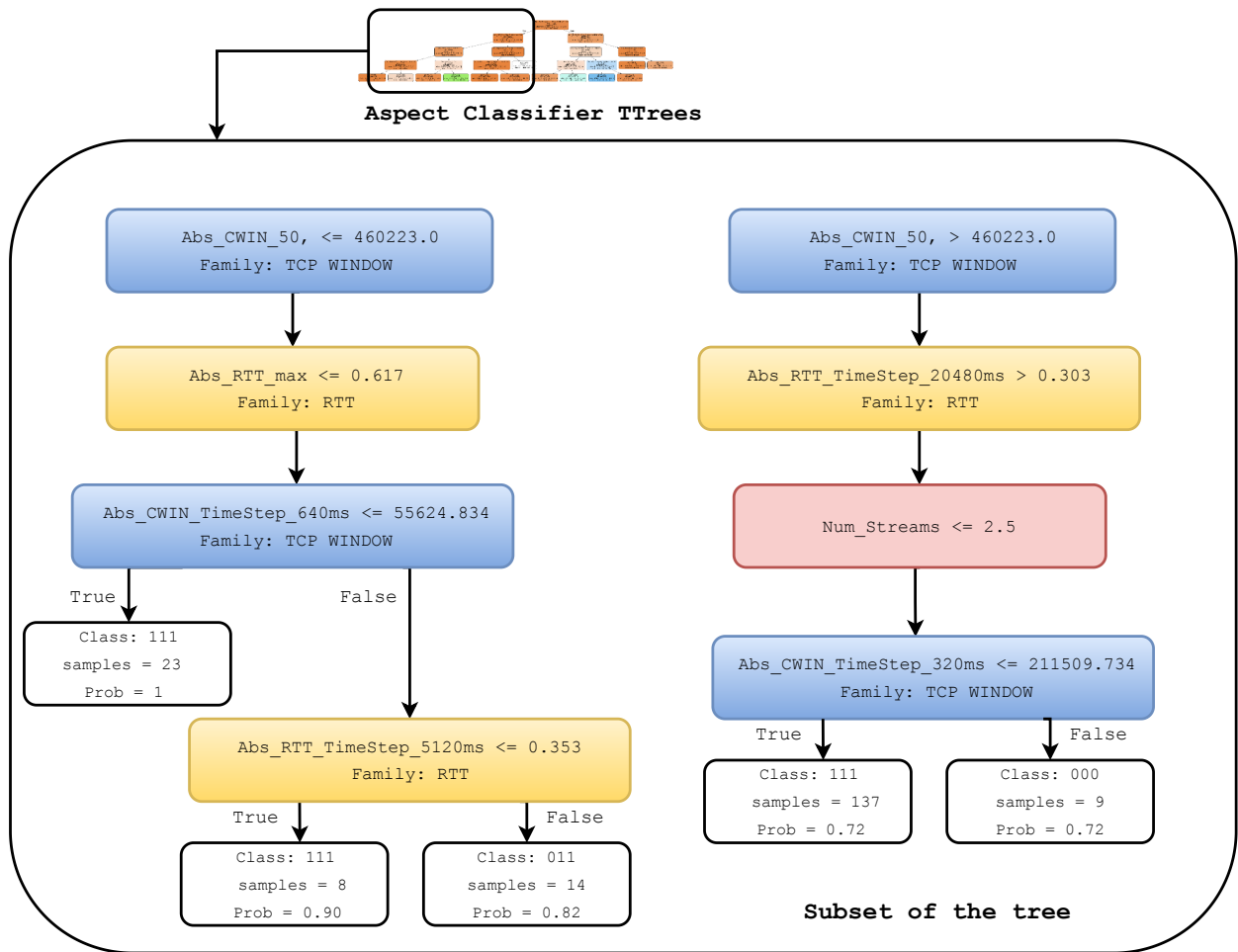


Figure 13: Zoom into a subset of the branches of the aspect classification tree C_2 built by TTrees for Dataset #4. The figure also reports which family of features the branching variable belongs to (see Table 2)

689 selection was eventually advantageous in terms of accuracy, as visible in Table 3.

690 With Dataset #4, trees C_1 and C_2 are smaller than with the other datasets, which allows us to analyze them more
 691 in detail. A general analysis of the classification trees C_1 (partially depicted in Fig. 12) and C_2 (Fig. 13) shows that
 692 critical QUIC features were properly selected to classify the performance of tests. In particular congestion window size
 693 (specifically, its 75th percentile) is at the root node of C_1 . Indeed, given the characteristics of the tests performed using
 694 different congestion algorithms and assuming similar radio conditions, it is clear that the main features differentiating
 695 each type of test has to be the congestion window related statistics. Other features like the number of streams,
 696 buffering/delay (RTT), packet loss or segment sizes are also correctly classifying each test. For example, RTT values
 697 above 100 ms identify lower throughput samples due to capacity limitations (high buffering). A higher number of
 698 streams is also normally identifying better samples in the knowledge tree C_1 built by TTrees for Dataset #4.

699 The aspect classification tree C_2 tells that, with high probability, we either have an anomaly due to a low congestion
 700 window (maybe due to packet loss), or a combination of the three identified network aspects. For the first case we have
 701 two leaves of the tree, and in both leaves we can see that CWIN and RTT in the initial five seconds were OK, but at
 702 the end of the test they were not. This might be due to packet loss. In the case of combined anomalies, the tree shows
 703 that there was low radio capacity or the congestion algorithm was too aggressive. In both cases we have identified
 704 the cause of anomalous performance of QUIC with the same methodology used to identify anomalous performance
 705 of TCP.

Table 6: A summary of the knowledge tree’s (C_1) performance in TTrees utilizing Facebook & Google as a service in all the countries part of the MONROE project with all their operators.

Services	Country	Operators	Accuracy (%)	# Samples	# Anomalous samples
Facebook (Dataset #5)	Sweden	<i>op0_sw</i>	96.0	183240	3162
		<i>op2_sw</i>	96.6	235568	3346
	Norway	<i>op0_no</i>	96.5	160290	2235
		<i>op1_no</i>	96.6	130334	2139
		<i>op2_no</i>	95.4	128804	2688
	Italy	<i>op0_it</i>	96.9	66420	1008
		<i>op1_it</i>	96.5	83966	1049
	Spain	<i>op0_es</i>	97.5	26795	320
		<i>op1_es</i>	95.5	30575	380
		<i>op2_es</i>	96.3	86450	1044
Google (Dataset #6)	Sweden	<i>op0_sw</i>	92.7	146327	3061
		<i>op2_sw</i>	98.9	180784	953
	Norway	<i>op1_no</i>	99.2	127628	469
	Italy	<i>op0_it</i>	96.8	66420	1008
		<i>op1_it</i>	96.5	83966	1049
	Spain	<i>op0_es</i>	98.4	24530	117
		<i>op1_es</i>	91.4	2078	32718
		<i>op2_es</i>	88.5	77729	7145

5.6. Modifications needed to deal with imbalancedness of datasets

To further extend the evaluation of TTrees, we employ Datasets #5 and #6, collected with MONROE. These datasets are collected from real operational mobile networks over which we do not have control. However, we have tested the connectivity to Facebook and Google webpages, which are highly replicated across the network, and easily reachable without having to traverse large portions of the cellular and core networks of the mobile operators to which MONROE probes were connected. As such, these datasets suffer imbalancedness, meaning that TCP problems are rare, while radio conditions might impair performance with much higher frequency, especially for mobile MONROE probes. Therefore, the analysis of Datasets #5 and #6 is suitable to showcase possible drawbacks due to the use of simple ML mechanisms in TTrees, and in particular the clustering operated with k -means, which is known to suffer in case of parsing imbalanced datasets [49, 50].

Table 6 summarizes the results obtained with TTrees in terms of accuracy and the number of samples concerning the knowledge tree for Datasets #5 and #6. The table shows results classified according to the origin of the data collected, i.e., per country of test and per operator. Similarly, Table 7 showcases the performance of TTrees when it comes to the aspect classification tree for the same datasets. For both datasets, the accuracy obtained with TTrees in both trees C_1 and C_2 is high, and the number of anomalies is as limited as it could be expected. However, Table 8 reveals that TTrees does not always manage to reach a fully meaningful classification of anomaly causes. In particular, the fourth and fifth columns of Table 8 show that the number of aspects identified with TTrees, and of the corresponding families, by means of k -means clustering is small for these datasets. This is a problem because the less aspects and families are extracted by the algorithm, the poorer is the precision with which the cause of an anomaly can be explained. Moreover, Table 8 does not report the results for some of the tested networks because in those cases there were not enough aspects to generate C_2 . The problem with the analysis of these datasets is that the anomalies found by means of C_1 are strongly imbalanced. For instance, Figure 15 shows that k -means finds most of the anomalies of C_1 concentrated in one of the two clusters extracted after applying binary classification using the most relevant features. In some scenarios, such as the ones displayed in those figures, this led to the impossibility of finding a set

Table 7: A summary of the aspect classification tree’s (C_2) performance with TTrees utilizing Facebook & Google as a service in all the countries part of the MONROE project with all their operators.

Services	Country	Operators	Accuracy (%)	# Samples	# Anomalous samples
Facebook (Dataset #5)	Sweden	<i>op0_sw</i>	98.5	183240	104
		<i>op2_sw</i>	98.9	235568	295
	Norway	<i>op0_no</i>	98.9	160290	2235
		<i>op1_no</i>	98.6	130334	2139
		<i>op2_no</i>	98.8	128804	770
	Italy	<i>op1_it</i>	98.7	83966	5
	Spain	<i>op0_es</i>	99.0	26795	18
		<i>op1_es</i>	98.9	30575	38
<i>op2_es</i>		99.0	86450	1044	
Google (Dataset #6)	Sweden	<i>op0_sw</i>	98.5	146327	105
		<i>op2_sw</i>	98.9	180784	322
	Norway	<i>op1_no</i>	99.6	127628	23
	Italy	<i>op0_it</i>	98.9	66420	89
		<i>op1_it</i>	99.0	83966	0
	Spain	<i>op0_es</i>	99.5	24530	0
		<i>op2_es</i>	95.5	77729	350

of features that fell under the balancedness threshold, and, thus, no aspects were obtained for labelling anomalous classes in C_2 . Therefore, applying k -means on these imbalanced C_1 anomalies can incur performance problems, as actually witnessed in our experiments.

Note that we have designed our CIAN methodology in a modular way, and the TTrees implementation relying on k -means is only a possible implementation with some strong advantages (it relies on simple algorithms and is interpretable) and some cons due to the limitations of the selected algorithms. However, in CIAN, methodological phases are well separated and the overall system is highly modular. Therefore, tools and algorithm used for a particular implementation can be easily exchanged guaranteeing no inter-dependency between technologies at different steps. Considering this, we can think of combating imbalancedness through a relatively simple modification of TTrees to implement CIAN. The alternative implementation of the methodology can simply differ from TTrees in how problems are clustered. In particular, we tested a modified version of TTrees obtained by interchanging k -means with Gaussian Mixture Model (GMM) [51], which is robust to imbalancedness and operates clustering in a very flexible way, although with less interpretable results.

5.6.1. Modified TTrees with GMM clustering

It is important to bear in mind that GMM is a model-based clustering approach that works with the concept of soft clustering, meaning that data is not just assigned to a cluster, but instead it has a likelihood of belonging to it [51]. Therefore, the objective of applying GMM to implement CIAN is to maximize this probability, a concept which can be directly associated to that of cluster tightness measured through inertia. Indeed, the previously defined TTrees implementation computes the inertia score for each pair of attributes. Instead, GMM requires to use of a different metric to draw performance comparisons between different pairs of attributes. The most well-known score for the selection of the optimal model in model-based clustering is the Bayesian Information Criterion (BIC) [52], so we use the BIC metric. Like inertia, the BIC metric is a penalty measure, which has to be minimized.

To show that GMM and k -means behave in a substantially different way, we evaluate the behavior of BIC and inertia obtained by applying GMM and k -means to form C_2 with the data relative to two operators within what collected in Dataset #5. Specifically, Figure 14 depicts the clustering performance values for the pairs of attributes extracted for Dataset #5 with *op0* in Spain, using GMM or k -means, respectively. Given that having lower inertia and

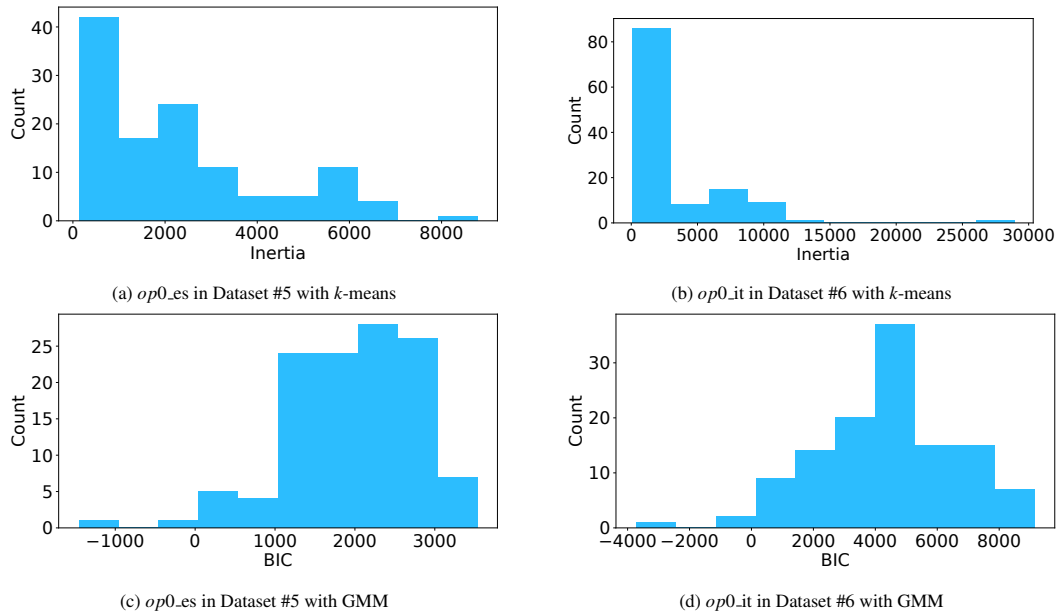


Figure 14: Histogram of Inertia (a, b) and BIC (c, d) values extracted during aspect clustering for all relevant aspects and with two operators in Spain and Italy. The lower the BIC or inertia the cleaner cluster grouping and splitting of samples is overall

756 BIC values translates into better separation between clusters, we see that the BIC metric reports values that are clearly
 757 lower than the ones extracted with inertia, and, more importantly, grouped within a considerably smaller range. This
 758 difference in the consistency of clustering performance may be an indicative of unstable behaviour of *k*-means, thus
 759 showing the presence of imbalancedness in the dataset. Moreover, according to the analysis conducted, low inertia
 760 would correspond to values between a 40 and 100, which as it is shown in the image, it is not the case for most
 761 attributes. This translates in the obtained clusters not being the most optimal, highlighting an under-performance of
 762 *k*-means compared to previous datasets and resulting of this imbalancedness. Figure 14 shows that the same analysis
 763 carried out for another operator and dataset (in Italy and with Dataset #6) yields an upper threshold for inertia values
 764 higher than the one obtained with BIC and a similar tendency in the distribution of data compared with the scenario
 765 mentioned previously. This proves that the imbalancedness of the dataset is consistent and observable across multiple
 766 operators and scenarios and thus it is a direct result of the nature of the experiments conducted.

767 Thus, using BIC, GMM is able to identify substantially different behaviors across the data collected for one
 768 operator or the other, while *k*-means might not be able to do so, at least when the datasets are imbalanced. Figure 15
 769 clearly shows that the anomalies obtained from C_1 are indeed imbalanced for the two cases analyzed in Figure 14.
 770 During the clustering phase of the methodology, samples classified as anomalous are fit into one of two clusters for
 771 each aspect: a cluster of problematic samples and one of compliant (non-problematic) samples. Since “N Class 0”,
 772 which is the ratio of samples falling under the non-problematic class for this binary classification problem, only takes
 773 values that are very close to 0 or 1, the anomalies are indeed imbalanced, as opposed to being evenly distributed.

774 By comparing aspect classification performance across more scenarios, the sixth and seventh columns of Table 8
 775 show that using GMM instead of *k*-means, with Datasets #5 and #6 finds more aspects. The last two column show
 776 that GMM finds most of the aspects found by *k*-means and identifies additional ones. This shows that GMM tends to
 777 be more general than *k*-means. However, the results obtained with *k*-means are more interpretable. Interpretability in
 778 decision trees is tightly correlated to the tree’s size. Thus, the higher the number of nodes, the less interpretable the
 779 tree is. The significant difference in the number of aspects identified with GMM translates in an increase in the depth
 780 and branching of C_2 . Therefore, although GMM was able to compute a solution in those particular cases, it yielded
 781 clusters and an aspect classification trees harder to traverse and understand. This relates to the fact that although GMM
 782 can handle complex patterns in data, it does so at the expense of interpretability [53].

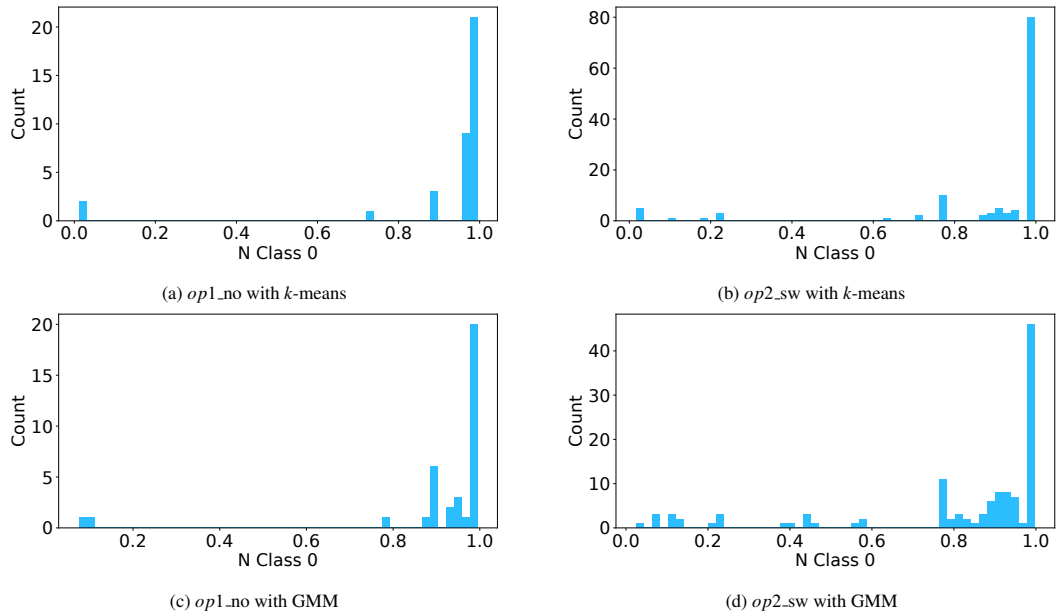


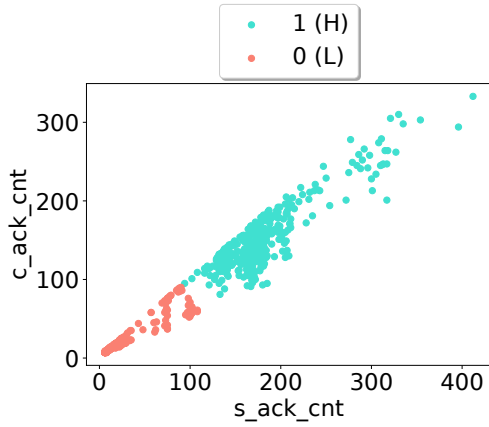
Figure 15: Histogram of the ratio of samples falling under the non-problematic class during aspect clustering for all relevant aspects in Dataset #5 and with two operators in Sweden and Norway. As it can be seen, most aspects, especially with k -means, yield imbalanced clusters of anomalies.

5.6.2. Take away message

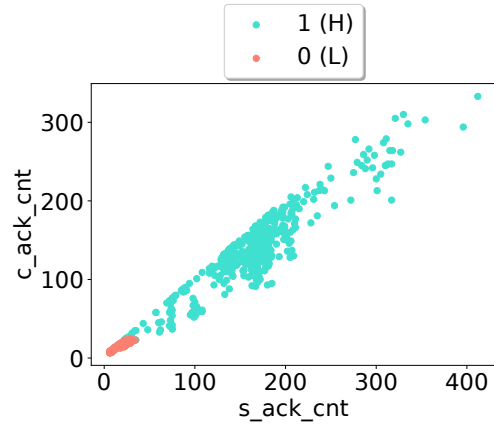
In case the dataset is balanced k -means and GMM convey similar clustering each grouping the clusters based on how the algorithm works underneath as shown in Figures 16 and 17. This is true for Datasets #1 to #4. However, Table 8 tells that the differences in the aspect families extracted from the application of k -means and GMM are significant in Datasets #5 and #6. The table reports the number of intersecting aspects and the ratio between intersecting aspects and the size of the smallest aspect set among the ones identified by k -means and GMM. The value of the ratio is, in most cases, close to one, which hints to the fact that the intersection of aspects found with the two approaches is similar to the smallest set identified by either GMM or k -means. Moreover, it is significant that in all the subsets available in Datasets #5 and #6, GMM finds a group of attributes equal or larger to the one extracted with k -means. In practice, as shown in the table, all the aspect families detected by k -means are extended by GMM, which is indicative of the fact that k -means identifies a subset of the data extracted with GMM. Thus, differently from GMM, k -means enforces a more aggressive filtering behaviour when selecting aspects, and this behavior is strongly affected by the balancedness filter used in the proposed methodology.

This difference between using GMM or k -means is clearly displayed, e.g., in Figure 15, in which we can see how GMM is able to identify a significantly higher number of aspects closer to an ideally balanced distribution of samples. It is also important to note that in some cases, k -means does not find a sufficient number of aspects and is not capable of identifying the root for the existence of anomalies. In fact, when TTrees with k -means does not obtain an aspect label, no aspect families are highlighted in C_2 , and TTrees cannot identify clearly whether a group of anomalies as classified by the decision tree C_1 is problematic or not for the target KPI. Indeed, the fact that some operators are not included in Table 8 alludes to the failure of some experiments utilizing k -means. We observed that TTrees with the latter clustering algorithm was not able to generate C_2 due to a highly imbalanced number of anomalies obtained from C_1 . On the contrary, the application of GMM as the clustering algorithm in the same scenarios always resulted in a complete aspect classification tree. Note also that, although GMM was able to compute a solution in those particular cases, it yielded clusters and an aspect classification trees harder to traverse and understand. Therefore, at the end of the day, if imbalanced anomalies are encountered, a feasible approach would be to switch the clustering algorithm to one that is able to handle this particular data distribution, thus, allowing for a higher flexibility.

Ideally, we want an easy-to-interpret separation to distinguish between anomaly samples, and this is why in TTrees we have opted for k -means. However, in case of imbalanced datasets, a more complex separation of anomalies into

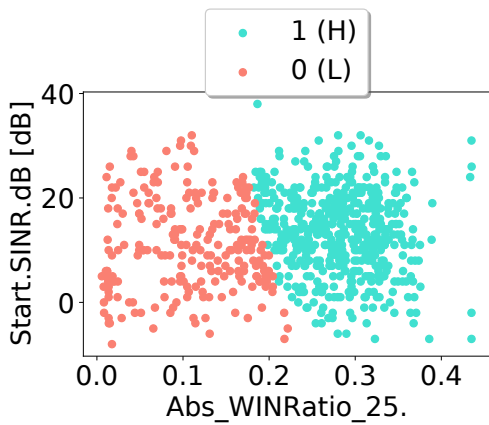


(a) Aspect 1 (Packet anomalies) with k -means

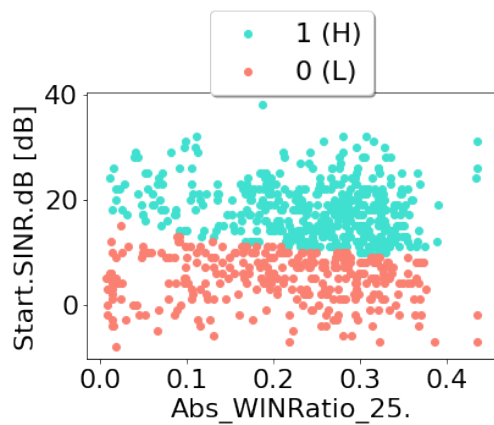


(b) Aspect 1 (Packet anomalies) with GMM

Figure 16: Cluster grouping using k -means and GMM as the cluster algorithm using Google Dataset and *op0.it*



(a) Aspect 1 (TCP Congestion Window Vs Signal-To-Noise-plus-Interference Ratio (SINR)) with k -means



(b) Aspect 1 (TCP Congestion Window Vs Signal-To-Noise-plus-Interference Ratio (SINR)) with GMM

Figure 17: Cluster grouping Comparison using k -means and GMM as the cluster algorithm for Dataset #1

Table 8: A comparison of the number of aspects and aspect families extracted from C_2 between k -means and GMM for the Google & Facebook datasets in all the countries part of the MONROE project. The table also shows the number of aspects that are common to both algorithms as well as the ratio, calculated from the division of the intersection by the size of the smallest set of aspects between k -means and GMM.

Service	Country	Operators	# k -means Aspects	k -means Aspect families	# GMM Aspects	GMM Aspect families	# Aspect intersection	Ratio
Facebook (Dataset #5)	Sweden	<i>op0_sw</i>	6	RTT, TCP Window, Packet anomaly	8	RTT, TCP, TCP Window, Packet anomaly	4	0.667
		<i>op2_sw</i>	4	RTT, Packet anomaly	10	RTT, TCP, TCP Window, Packet anomaly	4	1.0
	Norway	<i>op0_no</i>	2	-	11	RTT, TCP, TCP Window, Packet anomaly	2	1.0
		<i>op1_no</i>	2	-	10	RTT, TCP, TCP Window, Packet anomaly	1	0.5
		<i>op2_no</i>	5	RTT, TCP, TCP Window	10	RTT, TCP Window, Packet anomaly	4	0.8
	Italy	<i>op1_it</i>	8	RTT, TCP, TCP Window, Packet anomaly	10	RTT, TCP, TCP Window, Packet anomaly	4	0.5
	Spain	<i>op0_es</i>	4	RTT, TCP	7	RTT, TCP, TCP Window	3	0.75
		<i>op1_es</i>	6	RTT, TCP	7	RTT, TCP, TCP Window	5	0.833
		<i>op2_es</i>	2	-	7	TCP, TCP Window, Packet Anomaly	2	1.0
Google (Dataset #6)	Sweden	<i>op0_sw</i>	10	RTT, TCP, TCP Window	10	RTT, TCP, TCP Window, Packet anomaly	4	0.4
		<i>op2_sw</i>	8	RTT, TCP, TCP Window	9	RTT, TCP, TCP Window	7	0.875
	Norway	<i>op1_no</i>	5	RTT, TCP Window, Packet anomaly	9	RTT, TCP Window, Packet anomaly	5	1.0
	Italy	<i>op0_it</i>	9	RTT, TCP, TCP Window	8	RTT, TCP, TCP Window	7	0.875
		<i>op1_it</i>	10	RTT, TCP, Packet Anomaly	12	RTT, TCP, Packet Anomaly	8	0.8
	Spain	<i>op0_es</i>	1	-	1	-	1	1.0

811 clusters is needed, and can be achieved by using GMM, at the expenses of interpretability. Therefore, at the end of the
812 day, if anomalies are balanced, k -means is a suitable tool to be used to implement CIAN, whereas, should imbalanced
813 anomalies be encountered, a feasible approach would be to switch to GMM. In practice, we could extend TTrees by
814 automatically checking the presence of imbalanced anomalies, and then automatically select the appropriate clustering
815 algorithm between k -means and GMM, or another one that is able to handle the particular data/anomaly distribution
816 observed in the dataset. This would allow for a higher flexibility in the implementation of CIAN.

817 6. Conclusions

818 In this article, we have developed, validated and applied a methodology to identify the root causes of network
819 anomalies by means of interpretable ML algorithms. Our methodology, named CIAN, allows for automatic identifi-
820 cation of the networking aspects that cause anomalous behaviors. Such anomalies cannot be immediately and directly
821 explained by observing the network features. In fact, the potentially large number of network aspects, even in the pres-
822 ence of a limited number of features and samples per feature, makes it impossible to manually inspect and correlate.
823 With CIAN, we have purposely developed a radically novel methodology and implemented it with Python. Specifi-
824 cally, we have leveraged the key observation that if a network behavior cannot be modeled (or learned by ML), it is
825 a symptom of network anomaly. We therefore easily identify performance anomalies, after which we are able to use
826 unsupervised ML and Kolmogorov complexity-inspired tools to smartly search for the aspects that are more relevant
827 to explain where the anomaly is rooted. In particular we have introduced a novel metric, miscoding, for the evaluation
828 of redundancy and relevance of features. The final outcome of our proposed methodology is a set of easy-to-interpret
829 classification rules, that, in case of performance anomaly, allow for automatically alerting the appropriate departments
830 for corrective actions. To do so, CIAN only needs little volumes of samples for the features. CIAN is not specifi-
831 cally designed for a network protocol or service or dataset type, as we have validated in this work by applying its

832 implementation, TTrees, to the analysis of TCP and QUIC with data gathered with different granularity, formats and
833 richness of experimentally collected metrics. Indeed, TTrees allows to fully automatize network troubleshooting with
834 high accuracy and fast training, as shown in this article with the help of real data gathered from operational cellular
835 networks in various countries. We have also evaluated the impact of imbalanced anomalies on CIAN and verified that
836 TTrees might suffer performance issues in that case, which can however be obviated by simply re-placing (possibly
837 by means of an automated software selection process) k -means by GMM in the clustering of anomaly aspects. This
838 highlights the flexibility and modularity of the CIAN methodology.

839 Acknowledgments

840 The work was partially supported by the Regional Government of Madrid (CM), under grant EdgeData-CM
841 (P2018/TCS4499, cofunded by FSE & FEDER). The work was also partially supported by the Spanish State Re-
842 search Agency (AEI) PID2019-109805RB-I00/ AEI/10.13039/501100011033 grant (ECID project).

843 References

- 844 [1] M. Moulay, R. G. Leiva, V. Mancuso, P. J. Rojo Maroni, A. Fernández Anta, Ttrees: Automated classification of causes of network anomalies
845 with little data, in: 2021 IEEE 22nd International Symposium on a World of Wireless, Mobile and Multimedia Networks (WoWMoM), 2021,
846 pp. 199–208. doi:10.1109/WoWMoM51794.2021.00033.
- 847 [2] P. Yang, Y. Xiao, M. Xiao, S. Li, 6G wireless communications: Vision and potential techniques, IEEE Network 33 (4) (2019) 70–75.
848 doi:10.1109/MNET.2019.1800418.
- 849 [3] J. P. Santos, R. Alheiro, L. Andrade, A. L. Valdivieso-Caraguay, L. I. Barona-López, M. A. Sotelo-Monge, L. J. Garcia-Villalba, W. Jiang,
850 H. Schotten, J. M. Alcaraz-Calero, Q. Wang, M. J. Barros, SELFNET framework self-healing capabilities for 5G mobile networks, Trans.
851 Emerg. Telecommun. Technol. 27 (9) (2016) 1225–1232. doi:10.1002/ett.3049.
- 852 [4] A. Asghar, H. Farooq, A. Imran, Self-healing in emerging cellular networks: Review, challenges, and research directions, IEEE Communica-
853 tions Surveys & Tutorials 20 (3) (2018) 1682–1709.
- 854 [5] O. Loyola-González, Black-box vs. white-box: Understanding their advantages and weaknesses from a practical point of view, IEEE Access
855 7 (2019) 154096–154113.
- 856 [6] A. Holzinger, M. Plass, K. Holzinger, G. Crisan, C. Pintea, V. Palade, A glass-box interactive machine learning approach for solving NP-hard
857 problems with the human-in-the-loop, arXiv preprint arXiv:1708.01104 (2017).
- 858 [7] S. Dhanorkar, C. T. Wolf, K. Qian, A. Xu, L. Popa, Y. Li, Who Needs to Know What, When?: Broadening the Explainable AI (XAI) Design
859 Space by Looking at Explanations Across the AI Lifecycle, in: Designing Interactive Systems Conference 2021, DIS '21, Association for
860 Computing Machinery, New York, NY, USA, 2021, p. 1591–1602. doi:10.1145/3461778.3462131.
861 URL <https://doi.org/10.1145/3461778.3462131>
- 862 [8] F. K. Došilović, M. Brčić, N. Hlupić, Explainable artificial intelligence: A survey, in: MIPRO, 2018.
- 863 [9] M. Langer, D. Oster, T. Speith, H. Hermanns, L. Kästner, E. Schmidt, A. Sesing, K. Baum, What do we want from Explainable Artificial
864 Intelligence (XAI)? – A stakeholder perspective on XAI and a conceptual model guiding interdisciplinary XAI research, Artificial Intelligence
865 296 (2021) 103473. doi:<https://doi.org/10.1016/j.artint.2021.103473>.
866 URL <https://www.sciencedirect.com/science/article/pii/S0004370221000242>
- 867 [10] T. Ogata, A. Takeuchi, S. Fukuda, T. Yamada, T. Ochi, K. Inoue, J. Ota, Characteristics of skilled and unskilled system engineers in trou-
868 bleshooting for network systems, IEEE Access 8 (2020) 80779–80791.
- 869 [11] N. Bostrom, E. Yudkowsky, The ethics of artificial intelligence, The Cambridge handbook of artificial intelligence 1 (2014) 316–334.
- 870 [12] O. Alay, A. Lutu, M. Peón-Quirós, V. Mancuso, T. Hirsch, K. Evensen, A. Hansen, S. Alfredsson, J. Karlsson, A. Brunstrom, A. Safari Kha-
871 touni, M. Mellia, M. A. Marsan, Experience: An open platform for experimentation with commercial mobile broadband networks, in: ACM
872 MobiCom, 2017. doi:10.1145/3117811.3117812.
- 873 [13] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas,
874 A. Passos, D. Cournapeau, M. Brucher, M. Perrot, E. Duchesnay, Scikit-learn: Machine learning in Python, Journal of Machine Learning
875 Research 12 (2011) 2825–2830.
- 876 [14] T. Marwala, Causality, Correlation and Artificial Intelligence for Rational Decision Making, WORLD SCIENTIFIC, 2015. arXiv:<https://www.worldscientific.com/doi/pdf/10.1142/9356>, doi:10.1142/9356.
877 URL <https://www.worldscientific.com/doi/abs/10.1142/9356>
- 878 [15] B. Huang, K. Zhang, J. Zhang, J. Ramsey, R. Sanchez-Romero, C. Glymour, B. Schölkopf, Causal discovery from heteroge-
879 neous/nonstationary data with independent changes (2019). doi:10.48550/ARXIV.1903.01672.
880 URL <https://arxiv.org/abs/1903.01672>
- 881 [16] A. McGovern, R. Lagerquist, D. J. Gagne, G. E. Jergensen, K. L. Elmore, C. R. Homeyer, T. Smith, Making the black box more transparent:
882 Understanding the physical implications of machine learning, Bulletin of the American Meteorological Society 100 (11) (2019) 2175 – 2199.
883 doi:10.1175/BAMS-D-18-0195.1.
884 URL <https://journals.ametsoc.org/view/journals/bams/100/11/bams-d-18-0195.1.xml>
- 885 [17] J. Pearl, The seven tools of causal inference, with reflections on machine learning, Commun. ACM 62 (3) (2019) 54–60. doi:10.1145/
886 3241036.
887 URL <https://doi.org/10.1145/3241036>
- 888

- 889 [18] D. Garreau, U. von Luxburg, Explaining the explainer: A first theoretical analysis of lime, in: Proceedings of the 23rd International Confer-
890 ence on Artificial Intelligence and Statistics (AISTATS), Vol. 108 of Proceedings of Machine Learning Research, PMLR, 2020, pp. 1287–
891 1296.
892 URL <http://proceedings.mlr.press/v108/garreau20a.html>
- 893 [19] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, in: Proceedings of the 31st International Conference on
894 Neural Information Processing Systems, NIPS'17, Curran Associates Inc., Red Hook, NY, USA, 2017, p. 4768–4777.
- 895 [20] A. Shrikumar, P. Greenside, A. Kundaje, Learning important features through propagating activation differences, in: Proceedings of the 34th
896 International Conference on Machine Learning - Volume 70, ICML'17, JMLR.org, 2017, p. 3145–3153.
- 897 [21] A. Binder, G. Montavon, S. Bach, K. Müller, W. Samek, Layer-wise relevance propagation for neural networks with local renormalization
898 layers, CoRR abs/1604.00825 (2016). [arXiv:1604.00825](https://arxiv.org/abs/1604.00825).
899 URL <http://arxiv.org/abs/1604.00825>
- 900 [22] R. Barco, L. Nielsen, R. Guerrero, G. Hylander, S. Patel, Automated troubleshooting of a mobile communication network using bayesian
901 networks, in: MWCN, 2002.
- 902 [23] R. M. Khanafar, B. Solana, J. Triola, R. Barco, L. Moltsen, Z. Altman, P. Lazaro, Automated diagnosis for UMTS networks using Bayesian
903 network approach, IEEE Transactions on Vehicular Technology 57 (4) (2008) 2451–2461. doi:10.1109/TVT.2007.912610.
- 904 [24] S. Rezaei, H. Radmanesh, P. Alavizadeh, H. Nikoofar, F. Lahouti, Automatic fault detection and diagnosis in cellular networks using opera-
905 tions support systems data, in: IEEE/IFIP NOMS, 2016.
- 906 [25] E. J. Khatib, R. Barco, A. Gómez-Andrades, I. Serrano, Diagnosis based on genetic fuzzy algorithms for LTE self-healing, IEEE Trans. on
907 Vehicular Technology 65 (3) (2016) 1639–1651. doi:10.1109/TVT.2015.2414296.
- 908 [26] G. Ciocarlie, U. Lindqvist, K. Nitz, S. Nováczki, H. Sanneck, On the feasibility of deploying cell anomaly detection in operational cellular
909 networks, in: IEEE/IFIP NOMS, 2014.
- 910 [27] Q. Liao, S. Stanczak, Network state awareness and proactive anomaly detection in self-organizing networks, in: IEEE Globecom, 2015.
911 doi:10.1109/GLOCOM.2015.7414141.
- 912 [28] A. Terra, R. Inam, S. Baskaran, P. Batista, I. Burdick, E. Fersman, Explainability Methods for Identifying Root-Cause of SLA Violation
913 Prediction in 5G Network, in: GLOBECOM 2020 - 2020 IEEE Global Communications Conference, 2020, pp. 1–7. doi:10.1109/
914 GLOBECOM42002.2020.9322496.
- 915 [29] S. Tang, J. Kong, B. Niu, Z. Zhu, Programmable Multilayer INT: An Enabler for AI-Assisted Network Automation, IEEE Communications
916 Magazine 58 (1) (2020) 26–32. doi:10.1109/MCOM.001.1900365.
- 917 [30] Y. Wei, M. Peng, Y. Liu, Intent-based networks for 6G: Insights and challenges, Digital Communications and Networks 6 (3) (2020) 270–280.
918 doi:<https://doi.org/10.1016/j.dcan.2020.07.001>.
919 URL <https://www.sciencedirect.com/science/article/pii/S2352864820302418>
- 920 [31] P. D. Grünwald, The Minimum Description Length Principle (Adaptive Computation and Machine Learning), The MIT Press, 2007.
- 921 [32] C. S. Wallace, D. L. Dowe, Minimum Message Length and Kolmogorov Complexity, The Computer Journal 42 (4) (1999) 270–283.
- 922 [33] Y. Yang, G. I. Webb, Discretization for naive-bayes learning: managing discretization bias and variance, Machine learning 74 (1) (2009)
923 39–74.
- 924 [34] S. García, J. Luengo, J. A. Sáez, V. López, F. Herrera, A survey of discretization techniques: Taxonomy and empirical analysis in supervised
925 learning, IEEE Transactions on Knowledge and Data Engineering 25 (4) (2013) 734–750.
- 926 [35] D. V. Carvalho, E. M. Pereira, J. S. Cardoso, Machine learning interpretability: A survey on methods and metrics, Electronics 8 (8) (2019)
927 832.
- 928 [36] L. Breiman, J. Friedman, R. Olshen, C. Stone, Classification and regression trees., Kluwer Academic Publishers, New York (1984).
- 929 [37] T. Hastie, R. Tibshirani, J. Friedman, The elements of statistical learning: data mining, inference and prediction, Springer, 2009.
- 930 [38] W.-Y. Loh, Fifty years of classification and regression trees, International Statistical Review 82 (3) (2014) 329–348.
- 931 [39] H. Peng, F. Long, C. Ding, Feature selection based on mutual information: Criteria of max-dependency, max-relevance, and min-redundancy,
932 IEEE Trans. Pattern Anal. Mach. Intell. 27 (8) (2005) 1226–1238.
- 933 [40] G. Brown, A. Pocock, M.-J. Zhao, M. Luján, Conditional likelihood maximisation: A unifying framework for information theoretic feature
934 selection, Journal of Mach. Learn. Res. 13 (2012) 27–66.
- 935 [41] R. Cilibrasi, P. M. Vitányi, Clustering by compression, IEEE Trans. on Information theory 51 (4) (2005) 1523–1545.
- 936 [42] D. C. Hoaglin, F. Mosteller, J. W. T. (Editor), Understanding Robust and Exploratory Data Analysis, 1st Edition, Wiley-Interscience, 2000.
- 937 [43] D. Arthur, S. Vassilvitskii, K-means++: The advantages of careful seeding, in: ACM SODA, 2007.
- 938 [44] P. M. B. Vitányi, F. J. Balbach, R. L. Cilibrasi, M. Li, Normalized information distance, arXiv preprint arXiv:0809.2553 (2008). [arXiv:
939 0809.2553](https://arxiv.org/abs/0809.2553).
- 940 [45] V. Mancuso, M. Peón Quirós, C. Midoglu, M. Moulay, V. Comite, A. Lutu, O. Alay, S. Alfredsson, M. Rajiullah, A. Brunström, M. Mellia,
941 A. Safari Khatouni, T. Hirsch, Results from running an experiment as a service platform for mobile broadband networks in Europe, Computer
942 Communications 133 (2019) 89–101.
- 943 [46] A. Langley, A. Riddoch, A. Wilk, A. Vicente, C. Krasic, D. Zhang, F. Yang, F. Kouranov, I. Swett, J. Iyengar, J. Bailey, J. Dorfman,
944 J. Roskind, J. Kulik, P. Westin, R. Tennesi, R. Shade, R. Hamilton, V. Vasiliev, W.-T. Chang, Z. Shi, The QUIC transport protocol: Design
945 and internet-scale deployment, in: Proceedings ACM SIGCOMM, 2017.
- 946 [47] R. Marx, W. Lamotte, J. Reynders, K. Pittevels, P. Quax, Towards QUIC debuggability, in: Proceedings of ACM EPIQ, 2018.
- 947 [48] R. Marx, Main logging schema for qlog, Internet-Draft draft-marx-qlog-main-schema-02, Internet Engineering Task Force, Work in Progress
948 (Nov. 2020).
949 URL <https://datatracker.ietf.org/doc/html/draft-marx-qlog-main-schema-02>
- 950 [49] H. Xiong, J. Wu, J. Chen, K-means clustering versus validation measures: A data-distribution perspective, Systems, Man, and Cybernetics,
951 Part B: Cybernetics, IEEE Transactions on 39 (2009) 318 – 331. doi:10.1109/TSMCB.2008.2004559.
- 952 [50] S. Sieranoja, K-means properties on six clustering benchmark datasets, Applied Intelligence 48 (12 2018). doi:10.1007/
953 s10489-018-1238-7.

954 [51] C. M. Bishop, Mixture density networks, 1994.
955 [52] G. Schwarz, Estimating the dimension of a model, *The Annals of Statistics* 6 (2) (1978) 461–464.
956 URL <http://www.jstor.org/stable/2958889>
957 [53] C. M. Bishop, *Pattern Recognition and Machine Learning* (Information Science and Statistics), 1st Edition, Springer, 2007.
958 URL <http://www.amazon.com/Pattern-Recognition-Learning-Information-Statistics/dp/0387310738%3FSubscriptionId%3D13CT5CVB80YFWJEPWS02%26tag%3Dws%26linkCode%3Dxm2%26camp%3D2025%26creative%3D165953%26creativeASIN%3D0387310738>
959
960