

RESEARCH

Open Access



A method for identifying faulty cells using a classification tree-based UE diagnosis in LTE

P. Muñoz^{1*}, R. Barco¹, E. Cruz², A. Gómez-Andrades¹, E. J. Khatib¹ and N. Faour²

Abstract

The latest advances in wireless technologies have led to a proliferation of data mobile devices and services. As a consequence, mobile networks have experienced a significant increase in data traffic, while voice traffic has shown nearly no growth. It is therefore essential for operators to understand the data traffic behavior at the user level in order to ensure a good customer experience. In the radio access networks (RANs), traditional solutions based on cell-level measurements are not adequate to analyze performance of individual users. Instead, novel alternatives such as the use of call traces and the definition of new user-centric indicators will provide detailed and valuable information for each connection. One of the key measurements related to data services is the user throughput. In this work, the user throughput is adopted as the main attribute to conduct diagnosis in the RAN, which has typically been the bottleneck for data services. To that end, a binary classification tree is proposed to determine the root cause of poor throughput in user-level data sessions. Then, this information is aggregated at the cell level in order to provide effective diagnosis of degraded cells. In particular, a correlation-based analysis of the cell status is proposed in order to identify abnormal cell behaviors in an automatic way. Evaluation has been carried out with datasets from live cellular networks. Results show that the proposed diagnosis approach is an effective means to identify the main factors that limit the user throughput in network cells.

Keywords: Self-healing, Fault diagnosis, Long-Term Evolution, Correlation, Self-Organizing Networks

1 Introduction

During the last years, the wireless data services have become the dominant traffic source in cellular networks. Behind this, there is an expansion of new mobile applications and a rapid growth in the number of subscribers, both motivated by the advances in cellular communication technologies and the development of user-friendly smartphones. According to a large network vendor [1], global mobile data traffic grew 69% in 2014 while the average smartphone usage grew 45% in the same year. This enormous increase in data traffic has forced operators not only to invest large amounts of money in new infrastructure but also to reduce operational expenditures (OPEX) in order to maintain the levels of user satisfaction.

To produce significant cost savings, one of the adopted solutions by standardization bodies was the creation of

the Self-Organizing Networks (SONs) [2], which provide a new concept of network management where the maintenance and optimization tasks are carried out mostly in an automated way. Typically, technical experts in these fields have to deal with hundreds of traffic measurements and performance indicators every day [3, 4]. The vast diversity and quantity of these metrics makes the operational work very complex. Thus, the use of automated techniques for cellular traffic data analysis is essential to reduce human effort while expertise can be focused on new areas, bringing additional value to the operator [5].

Traditionally, mobile operators paid their attention in providing a good quality of the voice service, since it was the main offered service. To ensure this Quality-of-Service (QoS), troubleshooting experts mainly monitored the call blocking and dropping rates at the cell level to measure the levels of *accessibility* and *retainability*, respectively, in the network. However, with the explosion of Internet services, the QoS of multimedia and data applications is given by the data rates experienced by the users,

*Correspondence: pabloml@ic.uma.es

¹Communications Engineering Dept., University of Málaga, Málaga, Spain
Full list of author information is available at the end of the article

where *integrity* metrics such as throughput and latency are essential traffic measurements [4]. The problem of throughput performance indicators is that they are often difficult to interpret because of their dependence on many factors. In particular, there are some aspects beyond the typical variables related to the radio environment (e.g., distance to base station, cell loading, user speed, etc.) that should be considered. First, unlike in traditional voice services, the mobile network is only one segment of the end-to-end connection in an IP world. For example, a router in the IP cloud that suffers congestion may influence the user perceived data rate. Second, the recent radio access technologies (RATs) such as Long-Term-Evolution (LTE) have included a class-based QoS model as a mechanism to differentiate between services, establishing various levels of service to the users. Third, the traffic pattern of new data services clearly impacts throughput measurements. Due to the increasing popularity of web navigation, streaming video, social networking, file sharing, online gaming, and other data services, there are significant differences in traffic patterns [6]. As a consequence, operators are investing a large amount of money to investigate traffic modeling and classification through packet inspection in order to better understand the characteristic of today's cellular data traffic. In addition, sophisticated traffic data filtering, processing, and correlation with other network metrics are also important features to identify root causes of any detected anomaly and increase the reliability of the network [7, 8].

The increasing complexity of network infrastructure and services has also led operators to be interested in managing performance at the user level, instead of the cell- or network level, with the aim of maintaining their competitiveness levels. Today's solutions based on per-cell performance counters are insufficient to perform adequate root-cause analysis. For this reason, the standardization bodies have proposed the use of user-centric indicators and call traces to support the optimization and troubleshooting processes [9]. With the Minimization of Drive Tests (MDT) described in [10], the collection of traffic measurements can be done in an autonomous manner. In other words, each device that is active in the network reports measurements and signaling messages (i.e., call events) to the base station. Unlike traditional drive tests, MDT avoids the use of expensive measurement equipment and it does not require human effort. The information provided by call traces and MDT is not aggregated and reflects the performance at the user level. However, operators can process this information and use it to identify problems with greater accuracy at higher levels (e.g., the cell level).

This paper presents a novel method for determining problems in cells using information at the user level. Specifically, the input of the method is given by the

metrics contained in the call traces that provide specific information about data sessions. One of these metrics is the user throughput, which will act as the driver attribute to determine those data sessions experiencing bad QoS. This work focuses on the most common radio causes that may impact user throughput in a cellular network (i.e., congestion, lack of coverage, and interference). In case the user throughput is degraded, other radio measurements obtained from call traces will also be analyzed in order to identify the cause of bad QoS. To automate the analysis, a method based on a binary classification tree has been proposed. The adopted tree is generated from the analysis of real datasets and expert knowledge. Given the diagnosis of each data session, this information is aggregated for each network cell with the aim of creating a cell status or profile. Such a cell profile is then correlated to other cell profiles in order to identify faulty cells. The required calculations to classify the data and then compute the correlation values are carried out in an external server which communicates with the Operational Support System (OSS) in the network to obtain the network metrics. The study is carried out with different datasets from LTE networks where the proposed approach is applied to diagnose abnormal cells whose most users are experiencing poor user throughput.

Compared to other techniques, such as self-organizing maps [8], classification trees are an appropriate method to improve the effectiveness of diagnosis systems, especially when faults have a critical impact on network performance (i.e., various metrics are affected). Note that this kind of faults should be attended before any other. For example, a congestion situation is a critical problem that is reflected by a high number of connected users, but it may also be reflected by a high level of interference in the scenario. As a consequence, the diagnosis system could wrongly identify this situation as an interference problem. However, if the congestion situation is alleviated, the excessive interference level is also reduced. Thus, it is important for troubleshooting experts to have a clear map of which problems are prioritized by the diagnosis system. With classification trees, critical faults receiving higher priority should appear closer to the top of the tree. According to this, the classification trees enable fast visualization of the importance and prioritization of each fault.

The rest of the paper is organized as follows. The state-of-the-art is discussed in Section 2. The measurement setup and the real datasets to build and evaluate the proposed method are described in Section 3. The generation of the classification tree and the correlation-based method are covered in Section 4. The proposed method is evaluated using a real dataset in Section 5. Finally, Section 6 summarizes the main conclusions of this work.

2 Related work

The operational tasks in network management can be divided into three stages [11]: an initial measurement activity, a decision-making process, and lastly, a phase in which one or more actions are applied to the network. The first step means a continuous activity where a multitude of measurements are collected through different sources, including network counters and probes. Key Performance Indicators (KPIs) are continuously collected from network cells and then evaluated for optimization and troubleshooting purposes [12–14]. The main drawback of using those KPIs is that individual user performance may be lost if the data aggregation process at the cell level involves a considerable number of users. In addition, although they provide relevant information for managing the voice service, they are not enough for measuring the performance of data services. In this case, the use of user-centric KPIs would make the optimization and troubleshooting much more effective. For this reason, MDT is a feature introduced by the 3rd Generation Partnership Project (3GPP) allowing operators to utilize user devices to collect radio and traffic measurements in order to assess per-user level performance [15].

Some works addressing the MDT use case can be found in the literature. In [16, 17], MDT is utilized for coverage optimization, where the geo-localization of the user measurements enables powerful estimation and prediction of coverage holes. In [18], the location-aware radio measurements are employed for creating RF fingerprint databases which improves User Equipment (UE) positioning accuracy. In [19], a signal strength forecast method based on the classification and regression trees is proposed as another application of MDT. In [20], a KPI ranking system is proposed to significantly reduce the number of analyzed variables in MDT, while earlier work on QoS verification of MDT is described in [21]. In particular, the user experienced QoS in terms of throughput and its corresponding radio conditions are jointly analyzed. However, the evaluation is carried out by using a simulation tool instead of real traffic profiles and measurements.

The analysis of integrity performance (e.g., in terms of user throughput) in live networks has also gained attention in the research community. This kind of analysis can be carried out by means of field trials [22–24] or by using call traces (or MDT) [25]. However, in the case of field trials, the conclusions may not be representative of the real QoS experienced by the users. In the case of call traces, the correlation analysis presented in [25] is rather limited in terms of the number of radio measurements employed and no method for root-cause analysis was applied.

Due to the vast amount of data when per-user level information is collected, a new approach for network management is needed to address the requirements of the future fifth generation (5G). The main challenges

in the current SON paradigm to make 5G technically feasible has also been investigated in the literature. In [26–28], empowering SONs with Big Data techniques is studied with the aim of transforming big data into a readily useable knowledge base. In the field of self-healing, several research efforts have been devoted to the development of usable automatic detection and diagnosis systems [29]. On the one hand, various mathematical approaches have been applied to analyze network measurements, such as Bayesian networks [30, 31], Neural Networks [5, 8], Fuzzy Logic combined with Genetic Algorithms [32], linear prediction [33], correlation [34], and statistical analysis [35, 36]. However, these data-driven algorithms have been exclusively evaluated with per-cell level measurements, which may not be sufficient to manage the new data services.

The work in [8] whose aim is to diagnose problems at the cell level from traditional KPIs has been further investigated in [37] by employing call traces (as opposed to traditional KPIs) and applying a rule-based system to them. That work has been extended in [38], where a method based on Neural Networks (similar to that in [8]), is applied to diagnose the users. Then, from such a user diagnosis, a threshold-based method is applied to diagnose faulty cells. Due to its benefits, the present work also employs call traces as in [37, 38], instead of traditional KPIs. However, the proposed method improves the results presented in [38], whose method for cell diagnosis shows a clear dependence on thresholds settings. Unlike [38], the proposed method in this paper for cell diagnosis is based on the correlation to specific cell profiles so that the use of thresholds is avoided. In addition, the work in [37, 38] focuses on connections whose release has been abnormal, which can be considered a limitation in the case of data traffic, where the QoS evaluation during the session in both uplink and downlink is essential for diagnosis purposes.

On the other hand, there are some works in the literature that investigate the problem of cell outage by employing user-centric measurements [39–41]. Nevertheless, this information is commonly related to the radio environment (e.g., signal strength) derived from MDT functionality, while other measurements related to integrity performance such as user throughput are ignored.

Thus, there is a large fragmentation in references related to the abovementioned topics. For this reason, this paper attempts to unify such a fragmentation and overcome some limitations that have been found in previous works. First, the user throughput is used as a key indicator to estimate QoS of data traffic in mobile networks. Second, a method for troubleshooting network cells based on analysis of call traces is proposed. Third, rather than using traffic measurements from simulation tools or field tests, the evaluation is carried out

with large-scale datasets of real subscribers provided by operators.

3 Measurement setup and datasets

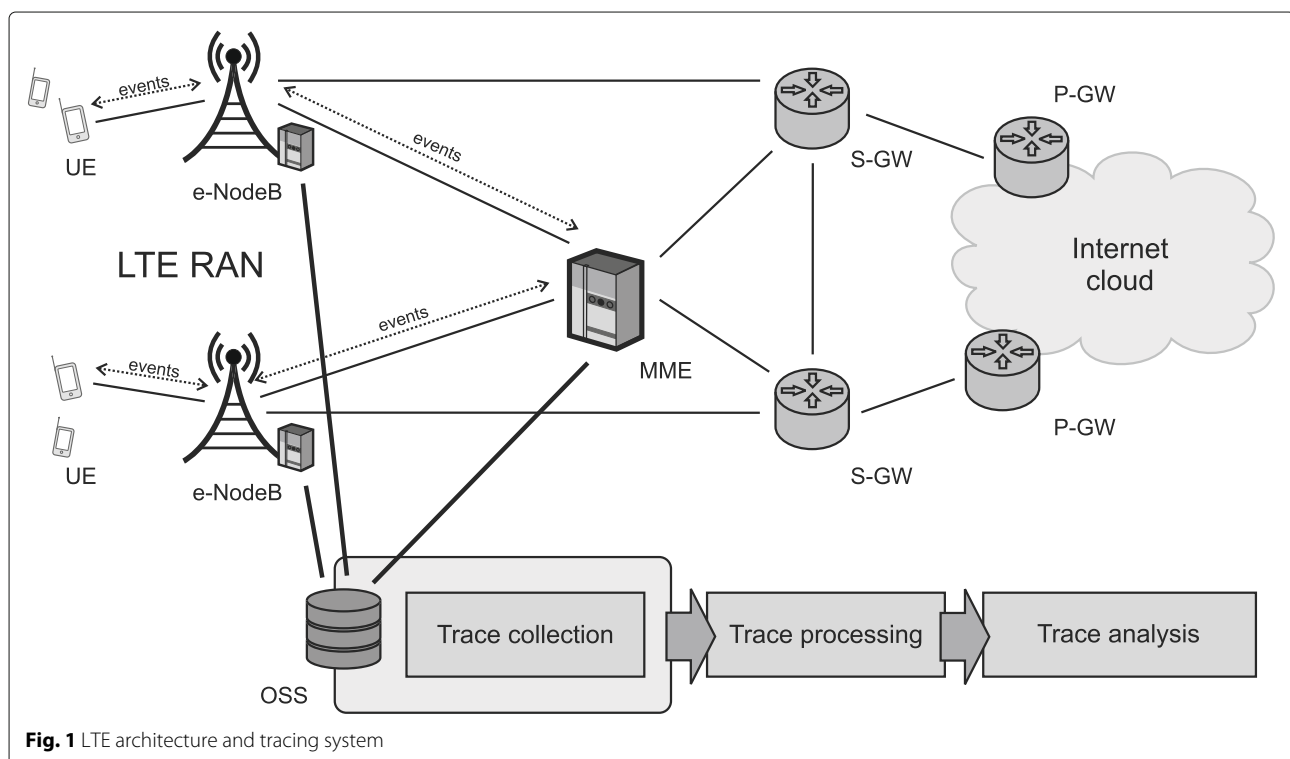
Figure 1 depicts the architectural model of 3GPP Radio Access Networks (RANs), together with the extensions to support the proposed data analysis framework. Since the data services are the dominant traffic source, this study focuses on LTE networks. The OSS is in charge of collecting events and measurements. These recordings are stored in files every report period, typically of 15 min duration. On the one hand, the connection between the OSS and the base station (e-NodeB) supports the collection of radio events. On the other hand, the connection between the OSS and the Mobility Management Entity (MME) enables the collection of user session details that can be used for mapping of radio events with measurements from another location of the network. Other network nodes (e.g., S-GW and P-GW) could also be employed to measure the quality of the services given to the user.

The tracing system provides detailed information at the session level for the UE. This information is contained in events which include measurements reported by the UE and signaling messages exchanged between the network elements. There are two different approaches for recording this information: UE traces and cell traces. The UE traces record events from those UEs which have been

selected by the operator. This allows operators to monitor a certain user if it is not getting the appropriate QoS. The cell traces record information from all UEs or a subset of UEs (provided some filters) in a selected cell. This approach can be used by operators for network optimization and troubleshooting purposes, since it provides larger datasets of per-user level statistics than UE traces. For this reason, cell traces have been used in this work.

After trace collection, the information has to be processed (trace processing in Fig. 1) in order to obtain relevant statistics at the user session level and KPIs with adequate granularity. For example, the data provided by an event that periodically reports a certain measurement can be aggregated in time to produce values in a longer time scale (e.g., at the user session level). A priori, the information could also be aggregated in the space domain (e.g., cells, cell clusters, or network). However, this may be equivalent to traditional approaches based on cell counters, where the measurements are typically aggregated at the cell level. Instead, the user-level statistics are directly analyzed by means of the trace analysis module. At this stage, the analysis of traces can be extended to higher levels, such as the cell level, keeping in mind that individual user performance should always be reflected to avoid hiding problems in the context of root-cause analysis.

The present work is carried out with large-scale measurements coming from performance recording applications which are executed in four different commercial



LTE network deployments. They are owned by different operators and located in metropolitan areas. The main characteristics of the networks are described in Table 1. Note that these values are calculated after filtering the data. Thus, the actual number of data sessions per cell is presumably higher; however, only the data sessions that provide reliable statistics have been considered in the study. The data filtering process will be explained in the next section. Another remarkable observation is that the cell bandwidth varies depending on the dataset. This may have an impact on user throughput. For example, it is observed that the average session throughput is higher in dataset 4 since the cell bandwidth in this network is greater.

4 Cell diagnosis based on per-user level traffic measurements

Cell traces provide very detailed information at the session level on every UE served by a cell. In particular, they give instantaneous values for a specific event. The events related to radio measurements play an important role in activities such as determining the root cause of malfunctioning UEs, analyzing dropped calls and optimizing resource usage and quality. In this work, the information contained in the traces is used to diagnose UEs whose performance in terms of throughput is degraded. This is carried out by means of a binary classification tree. The diagnosis made for each UE will serve to identify problems at the cell level. The following sections explain in detail how the classification tree is built and the subsequent cell diagnosis based on a correlation analysis.

4.1 Data input definition and filtering

In the context of data mining, the data pre-processing is essential to carry out troubleshooting and optimization tasks in an effective manner. Due to this, the first step in the construction of the classification tree is the definition of the input metrics from the information contained in the cell traces, followed by the filtering and cleaning processes. The metrics that have been defined in this work

are shown in Table 2. As observed, they represent the main aspects in a cellular network such as coverage, quality, and capacity in both directions of the radio link. By analyzing the distribution of these metrics over the network, most problems in the RAN can be diagnosed. For this reason, the number of selected metrics in this paper has been limited to 7. However, after building the classification tree, new metrics could be included in the tree in order to identify a larger number of problems, such as those related to mobility issues. To do this, the classification tree must comprise at least one leaf (i.e., a class) that represents the data sessions whose root cause remains unknown. From this leaf, the classification tree would be grown by introducing the conditions related to the new metrics. It is worth noting that the priority of the problems is affected by their position in the tree. The following paragraphs explain in detail each of these metrics as well as some aspects related to their filtering and cleaning processes.

To evaluate network performance, the 3GPP defines integrity as one of the basic categories for KPIs [4]. It attempts to measure how the RAN impacts the service quality provided to the users. Within this category, two different types of metrics are commonly defined: latency and throughput. In the downlink, the former is related to the delay experienced by the users, measured as the time from the reception of data in the e-NodeB to the transmission of the first packet over the radio interface. The latter is the data rate experienced by the users, measured as the data volume per elapsed time unit on the radio interface. In this work, the session throughput (Ses_Thp) is computed to identify users whose performance in terms of QoS is degraded. This requires the activation of the trace event that provides the corresponding data to compute the throughput for each UE in the downlink radio interface. In particular, this event periodically (every 2 s) reports the data rate for each radio bearer on the Downlink Shared Channel (DL-SCH) and Uplink Shared Channel (UL-SCH), which are the main transport channels for downlink and uplink data transfer, respectively. Let r_k be the k^{th}

Table 1 Characteristics of the datasets and network-related information

	Dataset 1	Dataset 2	Dataset 3	Dataset 4
Technology	LTE	LTE	LTE	LTE
Area	Urban	Urban	Urban	Urban
Cell type	Macrocells	Macrocells	Macrocells	Macrocells
Measurement period [min]	15	15	15	15
Cell bandwidth [MHz]	5 and 10	10	10	15
Number of cells	6421	238	318	655
Average number of data sessions per cell	4.07	4.07	4.38	5.19
Average session duration [s]	73.7	75.3	91.2	92.2
Average session throughput [Mbps]	5.1	5.3	6.7	11.6

Table 2 Selected metrics from cell traces

Metric	Link direction	Related magnitude
<i>Ses_Thp</i>	Downlink	QoS
<i>Num_UE_MHz_DL</i>	Downlink	Cell load per MHz
<i>RSRP</i>	Downlink	Signal strength
<i>RSRQ</i>	Downlink	Signal quality
<i>Num_UE_MHz_UL</i>	Uplink	Cell load per MHz
<i>Restr_Pwr_Ratio</i>	Uplink	Signal strength
<i>HARQ_Fail_Ratio</i>	Uplink	Signal quality

sample of data rate in the downlink collected during a session. The data rate r_k is internally computed by measuring the data volume in the so-called Time Transmission Intervals (TTIs).

Once the data rate samples at each measurement period k are gathered, *Ses_Thp* is calculated as an average of these samples. However, some of them should be excluded from the computation of this indicator in order to have a more accurate estimation of the throughput. The reason for this is that, with the new data services such as the online instant messaging, the level of bursty traffic has significantly increased. This obviously affects the evaluation of throughput performance indicators. On the one hand, there are a lot of TTIs in the traces where there are no data to transmit by the base station. These time intervals should be ignored in the computation of these indicators to make them independent of the file size. On the other hand, there are TTIs where the measured throughput is not representative of the data rate experienced by the UE. More specifically, given a traffic burst, the data volume transmitted in the TTI that empties the buffer (i.e., the last TTI in the burst) can negatively impact the average user throughput, especially if the amount of data in the last TTI is relatively much smaller than in the previous TTIs. The data volume in the last TTI also depends on the size of the packets at the user plane. As stated in [25], in the LTE downlink, almost 14% of data are transmitted in the last TTI, while around 40% of TTIs are last TTIs. To avoid the effect of the last TTI, in this paper, the data rate samples below a certain level, r_{min} , have been removed from the calculation of *Ses_Thp*. Assuming that N is the total number of samples in a data session, *Ses_Thp* can be formally expressed as:

$$Ses_Thp = \frac{\sum_{k=1}^N f(r_k > r_{min}) \cdot r_k}{\sum_{k=1}^N f(r_k > r_{min})}, \quad (1)$$

where $f(\cdot)$ is a function that returns “1” if the condition within the brackets is true, otherwise it returns “0.” In this paper, r_{min} is set to 250 Kbps.

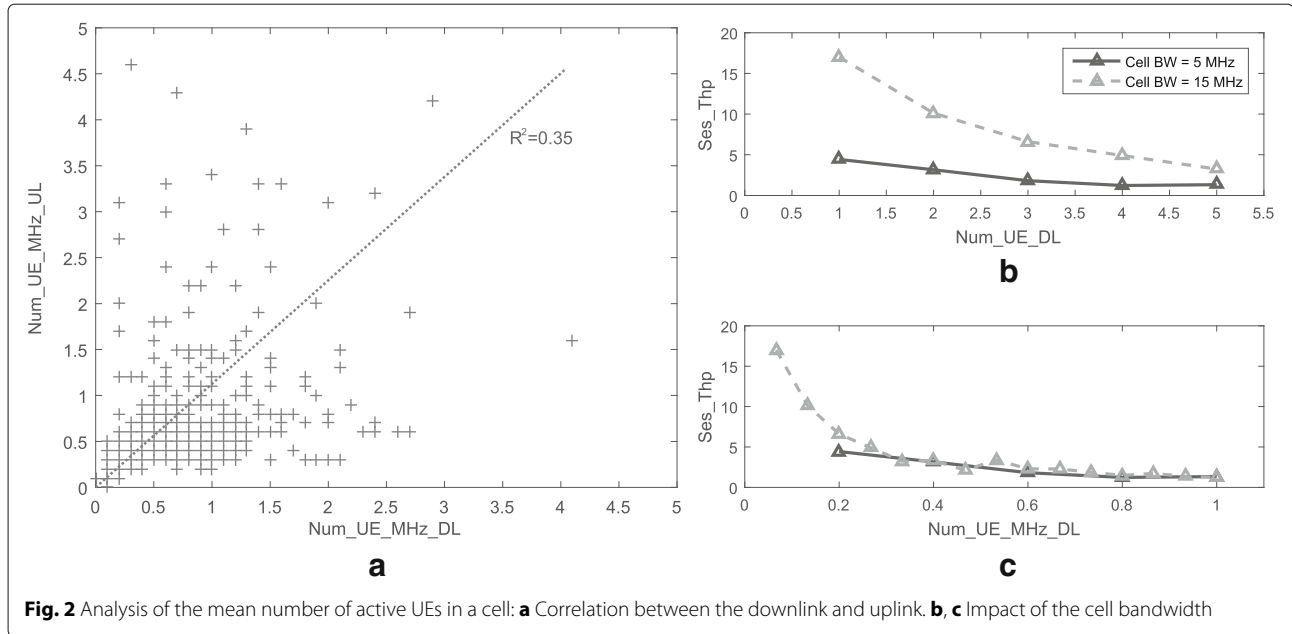
The metric *Ses_Thp* is also impacted by other factors, so that data filtering will be required at the session level.

For example, the Outer Loop Link Adaptation (OLLA) function [42] leads data sessions of short duration to abnormally low values of *Ses_Thp* [43]. The reason is that this mechanism adapts the modulation and coding scheme to provide certain block error rate in the connection. In particular, the OLLA starts with a conservative parameter setting and, after some iterations, the algorithm must converge. However, the target value is reached or not depending on the session duration. For short transmissions, the steady state of OLLA may not be reached and the modulation ramp-up is not long enough to compensate the initial setting of the algorithm. As a consequence, *Ses_Thp* will be negatively impacted. To avoid this problem, short data sessions should be removed from the dataset. Specifically, the event parameter available in the tracing system that measures the time to transmit a data burst excluding the last TTI can be used for this purpose [44]. Thus, those data sessions with this parameter below a certain threshold (set to 1500 ms in this work) are removed from the considered datasets.

While *Ses_Thp* is a measurement of the per-user QoS, the rest of metrics presented in Table 2 are measurements that describe the channel conditions in the radio environment for each UE. An estimation of the cell load in the downlink and uplink is given by *Num_UE_MHz_DL* and *Num_UE_MHz_UL*, respectively. They represent the mean number of UEs considered active in the corresponding link per TTI, calculated over one measurement period (2 s) and normalized by the cell bandwidth. Formally, it can be formulated as:

$$Num_UE_MHz_xL = \frac{Num_UE_xL}{Cell_BW}, \quad (2)$$

where x is “D” or “U” depending whether it is related to the downlink or the uplink, *Num_UE_xL* indicates the mean number of active UEs in the cell and *Cell_BW* expresses the cell bandwidth in MHz. The denominator can be obtained from the configuration management functional area, and the numerator is taken from an internal trace event in the cells. This event reports information about the cell traffic. It can be expected that a congested area would simultaneously lead to a high value of both downlink and uplink metrics (*Num_UE_DL* and *Num_UE_UL*). This effect is represented by the squared Pearson correlation coefficient (r^2) calculated for these two metrics, as shown in Fig. 2a. However, there are also frequent situations where only one link is congested, as illustrated in Fig. 2a, where the samples of dataset 1 (considering only cells with 10 MHz bandwidth) are represented. This means that the two metrics are meaningful from the diagnosis perspective. The metrics have also been normalized by the cell bandwidth since today’s LTE network deployments may have cells with different bandwidths, meaning that cells with higher bandwidth will provide higher



user throughput and support higher number of users. In Fig. 2b, c, the mean number of active UEs against the session throughput is represented in two cases: when the cell bandwidth is omitted, Fig. 2b, and when it is included in the definition, Fig. 2c. In the first case, given a number of active UEs, there is a big difference in *Ses_Thp* between cells with different bandwidth. Conversely, in the latter case, the values are very similar. For this reason, the cell bandwidth has been considered in the definition of cell load according to (2).

Another kind of metrics is related to the measurement of the signal strength. In the downlink, it is given by the Reference Signal Received Power (RSRP), defined as the average power of resource elements that carry cell specific reference signals over the entire bandwidth. RSRP levels typically range from around -75 dBm (close to an e-NodeB) to -120 dBm (at the cell edge). In the uplink, the information is taken from a trace event that reports the number of transport blocks on the Medium Access Control (MAC) level that are scheduled in the uplink, distinguishing between two cases. In particular, one is when the UE was considered to be power limited (counted by *Sched_TP_Restr_Pwr*) and another is when the UE was not limited in terms of power (counted by *Sched_TP_Unrestr_Pwr*). Based on these event parameters, the *Restr_Pwr_Ratio* is defined as follows:

$$\begin{aligned} \text{Restr_Pwr_Ratio} [\%] \\ = 100 \times \frac{\text{Sched_TP_Restr_Pwr}}{\text{Sched_TP_Unrestr_Pwr} + \text{Sched_TP_Restr_Pwr}}. \end{aligned} \quad (3)$$

A high value of this metric may represent a situation where the UE suffers from a lack of coverage since the UE is transmitting with the maximum power.

Finally, the measurement of the signal quality provides additional information of the radio environment, since it includes the interference component. In the downlink, it is given by the Reference Signal Received Quality (RSRQ), defined as:

$$\text{RSRQ} = L \times \frac{\text{RSRP}}{\text{RSSI}}, \quad (4)$$

where *RSSI* is the average total received power including the intra-cell power, interference, and noise, and *L* is the number of resource blocks over which the *RSSI* is measured (typically equal to the cell bandwidth). While the RSRP and RSRQ are reported by the UEs, the *RSSI* can simply be computed from RSRP and RSRQ. The range of RSRQ is normally from -19.5 to -3 dB.

In the uplink, an estimation of the signal quality is obtained from a trace event that reports the number of successful and unsuccessful transmissions at the hybrid automatic repeat request (HARQ) level in the uplink direction. The failed transmissions at the HARQ level are detected by means of a cyclic redundancy check (CRC). The provided information is also collected for each modulation format that is chosen by the UE. In this work, the statistics related to the quadrature phase-shift keying (QPSK) modulation have been considered since this modulation is typically used under the worst radio conditions.

Thus, from these event parameters, the following metric has been defined:

$$\begin{aligned} \text{HARQ_Fail_Ratio} [\%] \\ = 100 \times \frac{\text{HARQ_Fail_QPSK}}{\text{HARQ_Succ_QPSK} + \text{HARQ_Fail_QPSK}}. \end{aligned} \quad (5)$$

where *HARQ_Succ_QPSK* and *HARQ_Fail_QPSK* are the number of successful and failed transmissions, respectively, at the HARQ level in the uplink direction using the QPSK modulation.

To avoid certain dependence of performance evaluation on the datasets and, thus, facilitate the application of the proposed method to any LTE network, it is observed that the used metrics are not influenced by the network configuration. For example, the indicator to measure the traffic load is given per unit of bandwidth to avoid dependence on the system bandwidth. In addition, it is important that the number of considered cells (i.e., the network's size) in the datasets is large to include a significant number of faulty cells. Finally, note that the lack of data for any of the above described metrics in a few cells of a large network is a common situation, e.g., because the corresponding trace events have been erroneously deactivated. For this reason, as part of the data pre-processing, the cleaning function copes with incomplete and incorrect samples in the dataset. In this work, to avoid having inconsistent data, the entries (sessions) in the dataset with missed values are removed.

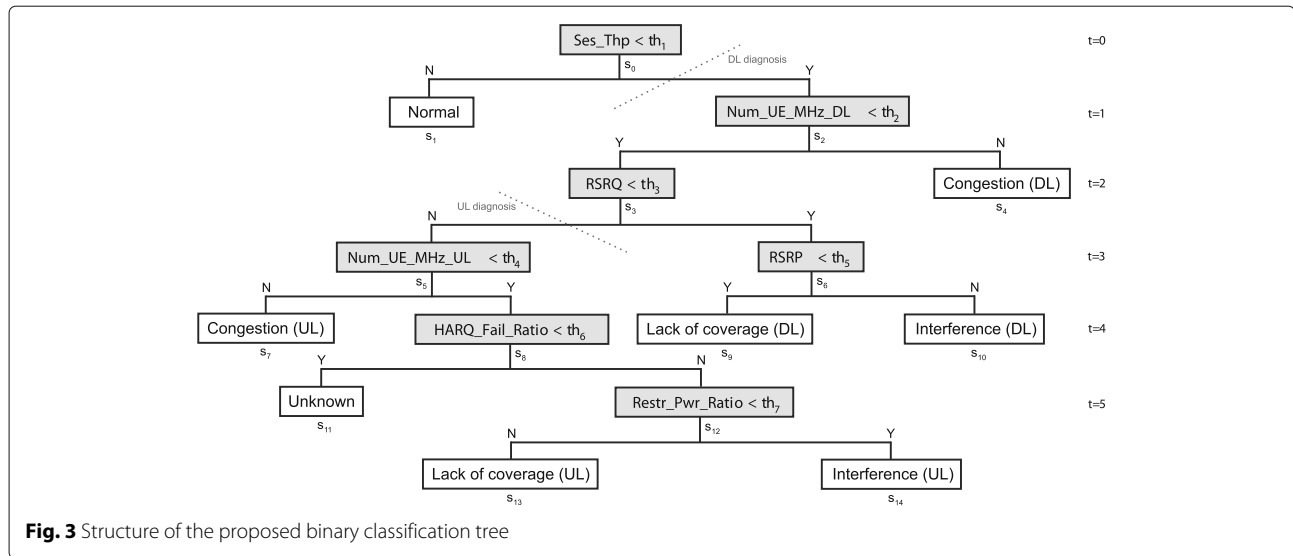
4.2 Construction of the classification tree

Traditionally, the RAN has been the bottleneck for data-hungry services. From the perspective of the QoS, it is important to identify if any degradation in the user throughput is due to the radio environment. The proposed classification tree aims at determining the root cause of abnormally low throughput values in user data sessions. This problem is addressed in the context of unsupervised learning since the training dataset is composed of unlabeled data, i.e., the root cause is a priori unknown. Considering a sufficient number of varied samples in a cellular network, the main problems related to radio aspects, such as congestion, lack of coverage and interference in either downlink or uplink, can be found in the given dataset. These root causes are determined by establishing adequate thresholds to radio KPIs. In the present section, the structure of the binary classification tree is proposed, while in the next section the adjustment of the thresholds is covered.

To better understand the binary classification problem, let Ω be the population of a training dataset, X_1, X_2, \dots, X_q be the set of q selected KPIs and C be the target attribute

that indicates the root cause and takes only a finite number of different values, i.e., $C = \{c_1, c_2, \dots, c_{q+1}\}$, where c_1 represents the normal behavior and q is the number of potential faults. Note that $q + 1$ categories can be distinguished using a binary classification tree with q KPIs. The goal of the tree is to generate a model $\phi(X_1, X_2, \dots, X_q)$ in the form of a decision tree for predicting the value of C from the values taken by the predictive variables. The solution is obtained by partitioning the Ω space into $q + 1$ disjoint sets, A_1, A_2, \dots, A_{q+1} . Thus, given a member of the space Ω denoted by ω , the predicted value of Y is c_j if ω belongs to A_j . The construction of the tree is based on iteratively splitting the dataset at each step by using one of the predictive attributes, X_1, X_2, \dots, X_q . In this work, a binary partition of the attributes is assumed since each attribute can be represented by two states, corresponding to normal and anomalous values. In most cases, one of the two children of the node will be a leaf node representing a specific problem. In addition, the nodes of the tree (referred to as s_j) obtained at each step t of the growing process define a partition that becomes finer with t . To make a binary partition, the logical operator "<" is used to compare each metric to a given threshold.

Figure 3 shows the logic of the proposed classification tree. Internal nodes and leaf nodes are colored in gray and white, respectively. The first partition of the tree, L_1 , acts as a detection step where data sessions with poor values of throughput are separated from those whose values are good. Thus, as illustrated in Fig. 3, the first leaf-node s_1 represents the group of data sessions with acceptable QoS, while the other branch of the split (i.e., node s_2) is subject to further analysis. As observed, the position of each variable in the tree will determine the priority of each problem in the diagnosis process. This establishes the way that troubleshooting experts operate when simultaneous problems happen in a network cell. To decide which variables are placed at higher levels in the tree, the statistical significance of the variables with respect to the primary variable, *Ses_Thp*, is first analyzed. More specifically, the predictive variables are used to fit a linear regression model where a set of coefficients are derived to describe the statistical relationship between *Ses_Thp* and the rest of KPIs. Then, the p values of the t statistics are computed for each coefficient in order to test the null hypothesis that the coefficient is equal to zero. In other words, a predictive variable X_j that has low p value is suitable to be included in the regression model since the changes in the variable are expected to be related to changes in *Ses_Thp*. On the contrary, a larger p value indicates that the changes in the predictive variables are not associated with changes in the primary variable. Table 3 shows the p values for each coefficient of the predictive variables in the linear regression model when 2000 samples from dataset 1 are used



as input. In the downlink, two metrics result in p values equal to zero, meaning that there exists a strong relation between the user throughput and such metrics. For this reason, in order to analyze the downlink user throughput (Ses_Thp), the radio conditions in the downlink are evaluated first in the proposed tree. This is represented by “DL diagnosis” in Fig. 3. After checking that the conditions in the downlink are good, the radio environment is then analyzed in the uplink (see “UL diagnosis” in the figure). For each link direction, the decision flow in the tree is started by checking the traffic load (nodes s_2 and s_5), because the p value obtained for $Num_UE_MHz_xL$ is zero in both link directions. The next conditions to be evaluated are those related to coverage and quality issues. In the downlink, it is observed that $RSRQ$ is notably more correlated to Ses_Thp than $RSRP$ since there is a big difference between their p values. As a consequence, the $RSRQ$ is used in the node s_3 to separate UEs with coverage and quality issues from those with good radio conditions in the downlink. The first group of UEs, related to node s_5 , are further diagnosed by analyzing the uplink through a decision flow that, for simplicity, is similar to the one applied in the downlink. The second group of UEs, related

to node s_6 , are separated by using the $RSRP$ with the aim of distinguishing UEs having a lack of coverage from those suffering from interference. Finally, in Fig. 3, there is one leaf node (s_{11}) representing the group of UEs that, having low values of Ses_Thp , are not associated to any of the considered root causes. This kind of UEs has been labeled as “Unknown” meaning that the root cause of their problem is not known.

Perturbation from external systems may produce false or very biased values of the metrics. However, this normally happens in a small number of UEs/cells compared to the size of the analyzed cell cluster/network. Thus, the “Normal” UE class and cell profile can be properly determined by the proposed method. If the external perturbation degrades the throughput but it does not produce any effect on the selected metrics, then the cause of this perturbation will be identified as “Unknown” by the system. On the contrary, if one or more selected metrics are degraded, then the system will probably produce diagnostic errors. Lastly, in case of a global impact of the perturbation, the expert team should be responsible for monitoring and troubleshooting this kind of problems.

The adjustment of each threshold th_j present in the classification tree is addressed in the next section. In particular, the thresholds from th_2 to th_7 are configured by calculating a certain percentile of data, given that degradation is associated to extreme values of the indicators. After this, the data distribution over the classes is analyzed for different values of th_1 . Based on this analysis, th_1 is configured.

4.3 Adjustment of the thresholds

The binary classification tree performs test on numeric features to divide the data into two groups: those whose

Table 3 The p -values of the predictive variables

Metric	p value	Related root cause	Link direction
$Num_UE_MHz_DL$	0.0	Congestion	Downlink
$RSRP$	$7.5e - 3$	Lack of coverage	Downlink
$RSRQ$	0.0	Interference	Downlink
$Num_UE_MHz_UL$	0.0	Congestion	Uplink
$Restr_Pwr_Ratio$	$1.7e - 12$	Lack of coverage	Uplink
$HARQ_Fail_Ratio$	$3.9e - 7$	Interference	Uplink

values for the variable are less than a threshold and those whose values are greater than or equal to the given threshold. Thus, the configuration of the thresholds determine the data distribution over the different rootcauses. The metrics placed in higher levels of the tree will have greater impact on the overall data distribution. In particular, the metric *Ses_Thp*, located at the top of the tree and used for detection purposes, affects the distribution of all the considered root causes. This variable, which is not related to a single problem, requires some practical experience to define the corresponding threshold. The problem of defining low/abnormal values of throughput is not straightforward, as it depends on many aspects such as the radio access technology or the user perception. Thus, a sensitivity analysis is needed to evaluate the overall impact on the diagnosis process. Conversely, the rest of variables in the tree are related to specific problems. To calculate the thresholds for such variables, an automatic technique referred to as Percentile-Based Discretization (PBD) [45] is used. In particular, each threshold is set to the X^{th} percentile of all the values of the corresponding metric in the dataset. Such a percentile represents an estimation of the percentage of data samples considered to be degraded over the total, assuming that the worst values of each variable most likely correspond to UEs having the associated problem. For example, UEs with bad RSRP are assumed to have a lack of coverage. The PBD technique is applicable when the training dataset is large enough to be statistically meaningful (as it is the case of the datasets used in this work), including not only the well-performing UEs but also the problematic UEs. Thus, the thresholds from th_2 to th_7 are adjusted by computing a certain percentile of the data.

Table 4 shows the corresponding percentiles for each metric of the datasets used in this work. Typically, if the metric is degraded by reaching high values, the 80th percentile has been selected. On the contrary, if degradation is given by the lower part of the range, the 20th percentile is taken. Note that the selected percentiles are not extreme values since there can be different levels of the magnitude of the problems. It is observed that the calculated percentiles, especially for the measurements RSRP and RSRQ, are similar to those used by troubleshooting

experts in practice. In addition, the obtained values are similar over the different datasets. Only dataset 4 reflects a notable deviation in the RSRQ values since its greater cell bandwidth has resulted in lower inter-cell interference conditions. Finally, note that the thresholds from th_2 to th_7 in the classification tree will be those derived from the PBD method, depending on the specific dataset. Thus, each dataset has its own set of thresholds to be used in the classification tree.

Once the mentioned thresholds have been adjusted, the sensitivity analysis of th_1 can be performed. Figure 4 shows the distribution in percentage of the data sessions over the root causes for each dataset with cell bandwidth equal to 10 MHz. The distributions are represented for three different values of th_1 . As observed, results between datasets 1 and 2 are quite similar in general. Compared to them, dataset 3 presents some differences in specific classes. In particular, the percentage of data sessions for the Normal class in dataset 3 is higher than in datasets 1 and 2, while the percentage of data sessions experiencing congestion in the downlink in dataset 3 is lower than in datasets 1 and 2. This observation reveals that the network in dataset 3 is less loaded and, as a consequence, better QoS is provided to the UEs. However, when $th_1=10$ Mbps is used, the differences between datasets in the percentage of Normal UEs are significantly reduced. This is mainly because of an increase of the percentage of UEs suffering from interference and UL congestion in dataset 3. Thus, setting th_1 to 10 Mbps would lead to similar distributions between normal and abnormal cases for the given datasets. The main disadvantage of using this setting is that the percentages of normal cases are significantly low, specifically below 20%, while the percentage of data sessions with unknown root cause is above 25%. Obviously, this represents a pessimistic view of the network status and may block the detection of those problems with major impact on QoS. On the contrary, the setting with the lowest tested value, $th_1=2$ Mbps, provides distributions where some fault classes are empty. This is also a bad choice provided that all the considered problems should be present in large datasets even in a small proportion. For these reasons, the setting $th_1=5$ Mbps has been selected in this work to provide more effective diagnosis of

Table 4 Defined thresholds based on the X^{th} percentile of the metrics for each dataset

Metric	X^{th} -ile	Dataset 1	Dataset 2	Dataset 3	Dataset 4	Units
<i>Num_UE_MHz_DL</i>	80	0.4	0.7	0.2	0.2	[UE/MHz]
<i>RSRP</i>	20	-118.0	-111.5	-112.1	-113.4	[dBm]
<i>RSRQ</i>	20	-18.0	-16.5	-17.5	-13.7	[dB]
<i>Num_UE_MHz_UL</i>	80	0.4	0.7	0.3	0.2	[UE/MHz]
<i>Restr_Pwr_Ratio</i>	80	99.0	98.7	99.9	99.8	[%]
<i>HARQ_Fail_Ratio</i>	80	10.4	9.2	16.0	13.4	[%]

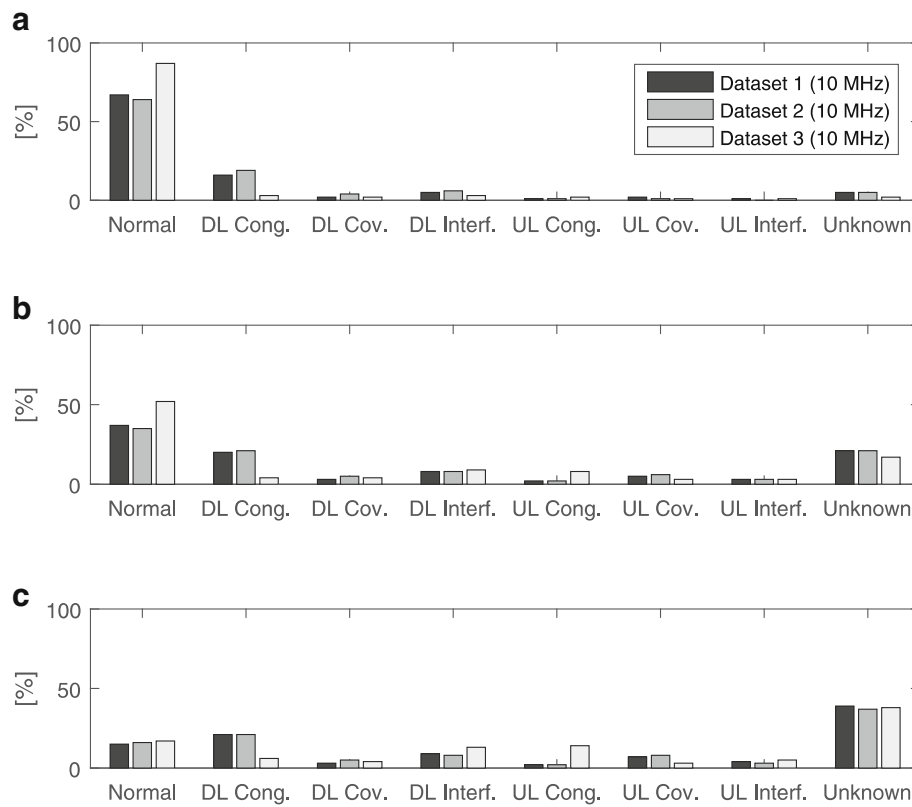


Fig. 4 Sensitivity analysis of th_1 for datasets with the same (10 MHz) cell bandwidth. **a** $th_1 = 2$ Mbps. **b** $th_1 = 5$ Mbps. **c** $th_1 = 10$ Mbps

the radio problems. In this sense, note that the Unknown class includes the UEs with bad QoS due to non-radio issues. Since most UEs are accessing to Internet services, the contribution to this class should be important. With the proposed threshold, the percentage of data sessions whose root cause is unknown remains relatively high but in a lesser extent than the Normal class, as expected.

Finally, Fig. 5 shows a comparison of the data distribution over the classes between the datasets with different values of cell bandwidth. The analysis is also made for each value of th_1 previously analyzed. As expected, there is a clear trade-off between the cell bandwidth (related to capacity) and the QoS of the data sessions. In particular, the higher the cell bandwidth, the higher the percentage of UEs belonging to the 'Normal' class is expected to be. On the one hand, it is noted that the main affected root cause due to a change in the cell bandwidth is congestion, while the data distribution over the rest of root causes is not significantly impacted. Thus, the proposed classification tree effectively establishes a relation between the cell bandwidth and congestion issues: a high contribution of the classes related to congestion suggests that there is a lack of extra bandwidth. On the other hand, the same tendency is observed regardless of the specific value of

the analyzed threshold, th_1 . Thus, the optimal setting for th_1 previously selected can be applied to scenarios with different cell bandwidth.

4.4 Cell-level correlation-based diagnosis

The final cause of the problem in a specific UE can be very diverse, especially when it connects to Internet. Abnormally low values of throughput can be registered, for example, when a router in the core network is malfunctioning, the server is overloaded or when an application in the UE is crashing. In this paper, the previous sections have proposed a practical system with the aim of determining the root cause of problematic data sessions when their radio conditions are not favorable. This means that isolated problems (such as the previously mentioned) are not considered in this work. On the contrary, the final goal of the proposed method is the diagnosis of problems in the radio environment that affect a significant number of UEs. More specifically, the diagnosis made for every UE is used to find localized problems in network cells. This is achieved by defining cell profiles (explained later), which represents the distribution of the diagnosed UEs over the predefined fault classes. Then, the comparison between different cell profiles using correlation techniques allows

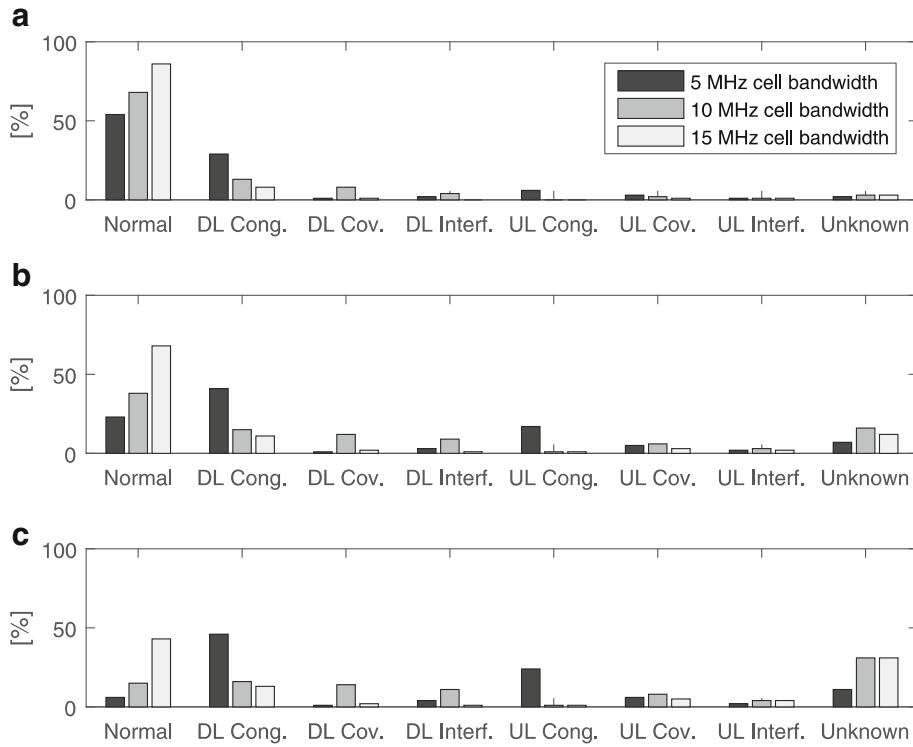


Fig. 5 Sensitivity analysis of th_1 for datasets with different cell bandwidth. **a** $th_1 = 2$ Mbps. **b** $th_1 = 5$ Mbps. **c** $th_1 = 10$ Mbps

identifying problematic cells. This differs from traditional approaches where individual measurement samples are aggregated in a non-selective way at the cell-level to provide the diagnosis. The main benefit of the proposed method is that the performance of the problematic UEs is not biased by the performance of the dominant UEs in a cell, since the diagnosis is first made at the user level. Essentially, the method is composed of five steps represented in Fig. 6.

The first step of the proposed method is to apply the binary classification tree to every UE in a certain cell in order to obtain the diagnosis at the user level. After this, the percentage of UEs having a specific issue (i.e., those UEs belonging to a certain class in the tree) is computed. In particular, the following cell-level indicator, denoted as ξ_j^k , is calculated for each class j in every cell k of the scenario:

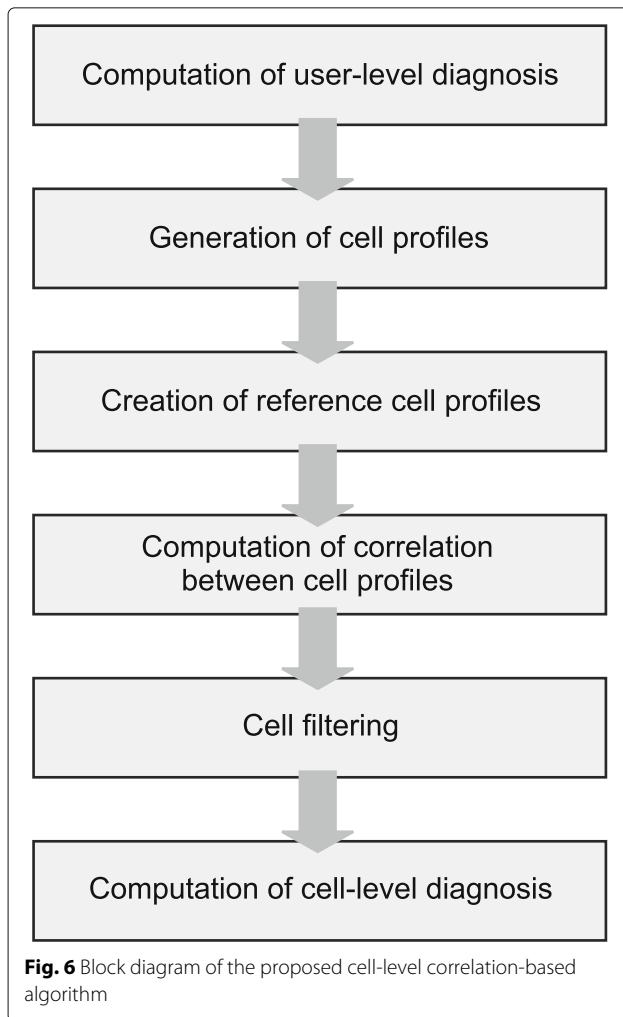
$$\xi_j^k [\%] = 100 \cdot \frac{|A_j^k|}{\sum_{i=1}^{q+1} |A_i^k|}, \quad \forall j = 1, 2, \dots, q+1. \quad (6)$$

where A_j^k is the set of data sessions associated to class j and served by cell k .

Once the above indicator is computed for the $q+1$ classes, the profile of a cell k is determined by the tuple $\xi^k = \{\xi_1^k, \xi_2^k, \dots, \xi_{q+1}^k\}$. A visual representation of the

cell profile provides an effective means of diagnosing cells. For example, in the case of cells without problems, the value of ξ is expected to be high for the Normal class and low for the rest of (fault) classes. However, this can be an overly simplistic view of the diagnosis process. Specifically, the profile of a healthy cell depends on many factors and it may vary from one network to another. To explain this, Fig. 4b can be seen as a representation of the profile of a network, in a similar way to cell profiles. Then, a close examination of this figure reveals that, although the UEs without problems are the dominant class in the network, there are still significant differences in the profiles between the different datasets. For example, the levels of interference in healthy cells of dataset 3 are expected to be higher than in the other two datasets. This does not necessarily mean that most cells of dataset 3 suffer from interference issues, rather it is merely an indication of the overall radio conditions in the network. For this reason, cell profiles should be adapted to every network in order to provide a better diagnosis.

Based on the above explanation, the next step of the method (Fig. 6) is the creation of the Reference Cell Profiles (RCPs). There exists one RCP that represents the faultless or normal behavior of a cell in the considered network, referred to as ξ^{normal} . Additionally, other RCPs



for faulty situations (e.g., $\xi^{\text{DL_Cong}}$, $\xi^{\text{DL_Cov}}$) can be generated. The main problem is that the RCPs must be created in an unsupervised manner because the actual performance of each cell (normal or faulty) is unknown. Clustering methods are useful to discover the most frequent cell profiles existing in the datasets. In this paper, data clustering is carried out by using the well-known k -means algorithm, which is the most popular clustering tool used in scientific and industrial applications. However, the RCPs can also be obtained by using other clustering methods, such as Self-Organizing Maps (explained later).

The k -means algorithm partitions the dataset into k clusters. Such a partition minimizes the sum of discrepancies between a point and its centroid expressed through appropriate distance (in this case, Squared Euclidean distance is used). Each cluster is represented by the mean (or weighted average) of its points, i.e., the so-called *centroid*. In this work, each centroid represents a RCP. The mechanism to choose the initial positions of the cluster centroids is based on randomly selecting k observations

from the dataset. Then, the algorithm iteratively executes three steps until convergence is reached: first, the distance of one data point to the centroids is calculated; second, the class label of the cluster centroid having the minimum distance is set to the given point; third, the centroids are updated based on the assignment. The convergence is reached when the partitioning error is not reduced by the relocation of the centroids or the maximum number of iterations allowed is reached. The number of clusters k establishes how many RCPs are identified. This number must be at least equal to the number of categories that have been defined in the classification tree for the user diagnosis (i.e., $q + 1$). A greater value of k can be used to identify new RCPs that would be associated by troubleshooting experts to one of the given categories. Note that an excessively high value might lead to wrong associations between the RCPs and the fault categories. For this reason, the value of k has been set to 12.

Table 5 shows the results of applying the k -means algorithm to dataset 1. Each row represents one obtained RCP or centroid. The first column is the identifier of the RCP, the second column shows the category assigned to the RCP, which can be the normal state or a fault, and the last columns provide the percentage of UEs associated to a specific issue for the given RCP. It is observed that RCPs from 1 to 8 are associated to categories where the dominant diagnosis at the user level corresponds to the labeled category at the cell-level. More specifically, RCP 1 is linked to the normal state, since not only the percentage of well-performing UEs is the highest, but also the percentage of UEs for the rest of (faulty) categories is very low. In line with the network profile represented in Fig. 4b, the second dominant category in RCP 1 is Unknown. This means that a cell normally has a tolerable percentage of abnormal UEs, labeled as Unknown, under normal cell conditions. RCPs from 2 to 8 are mainly related to specific problems, having a large percentage of UEs experiencing the dominant problem, but also a small percentage of UEs belonging to Normal and Unknown categories. Such particular distributions between categories, which varies among datasets, is an interesting reason to create specific RCPs to every network. Another reason comes from analyzing RCPs from 9 to 12, where it is observed that the Normal (and in many cases Unknown) category achieves similar levels to a certain fault category. Recall that the parameter k has been configured in order to find a reasonable number of these situations. This set of RCPs represents the most difficult faulty situations to be identified by the diagnosis system since the problematic UEs are not a clear dominant group within the coverage area of the affected cell. Moreover, in these cases, using traditional KPIs (i.e., cell-level measurements) may not be enough to identify the problem.

Table 5 Set of RCPs obtained by applying k -means with $k=12$

No. of RCP	Labeled categ. c_i	ξ_i [%]							
		Normal ($i=1$)	DL Cong ($i=2$)	DL Cov ($i=3$)	DL Int ($i=4$)	UL Cong ($i=5$)	UL Cov ($i=6$)	UL Int ($i=7$)	Unknown ($i=8$)
1	Normal	98.41	0.07	0.16	0.28	0.03	0.15	0.11	0.78
2	DL Cong	4.38	92.65	0.13	0.51	0.85	0.13	0.18	1.17
3	DL Cov	11.24	0.10	65.25	6.12	1.84	5.55	1.90	8.00
4	DL Int	0.80	0.00	0.85	93.38	0.66	0.16	0.71	3.44
5	UL Cong	3.59	3.24	1.21	3.71	85.23	0.98	0.51	1.54
6	UL Cov	3.66	0.00	0.78	2.03	1.06	86.83	1.22	4.43
7	UL Int	2.87	0.00	0.99	4.39	0.00	5.90	83.13	2.73
8	Unknown	1.77	0.03	0.25	0.67	0.08	0.35	0.32	96.54
9	DL Int	40.44	0.58	2.35	39.28	1.58	1.95	1.84	11.98
10	DL Int	7.96	2.48	4.51	17.30	1.71	11.13	8.35	46.57
11	UL Cov	53.69	6.05	2.30	1.73	8.19	15.40	7.77	4.87
12	Unknown	51.31	0.37	1.45	1.38	0.38	2.25	1.14	41.72

The next step of the proposed method is the calculation of a correlation coefficient that measures the similarity between a given cell profile and a reference cell profile. Specifically, the Pearson's correlation coefficient has been used in this work. Its possible values are in the range -1 to $+1$, indicating high linear correlation when it is close to any of those values. The plus sign represents positive correlation between the variables, while the minus sign indicates negative correlation. The Pearson's correlation between the cell profiles ξ^k and ξ^{ref} , referred to as $\rho_{k,ref}$, can be calculated according to this formula:

$$\rho_{k,ref} = \frac{\sum_{j=1}^{q+1} \xi_j^k \xi_j^{ref} - (q+1) \overline{\xi^k} \overline{\xi^{ref}}}{\sqrt{\left(\sum_{j=1}^{q+1} (\xi_j^k)^2 - (q+1) (\overline{\xi^k})^2\right) \left(\sum_{j=1}^{q+1} (\xi_j^{ref})^2 - (q+1) (\overline{\xi^{ref}})^2\right)}} \quad (7)$$

where $\overline{\xi^k}$ (and similarly $\overline{\xi^{ref}}$) is computed as:

$$\overline{\xi^k} = \frac{1}{q+1} \sum_{j=1}^{q+1} \xi_j^k. \quad (8)$$

To improve the accuracy of the proposed system, the cells are filtered before providing the final diagnosis. The filtering module is composed of two filters. First, the cells are filtered by the number of data sessions that have been used to compute their cell profile. Note that the cells with a low number of data sessions have been used to build the RCPs. Since this group of cells can be large in the dataset, they can provide (on average) meaningful information to discover frequent patterns. However, these cells are not adequate to be individually diagnosed since their

cell profiles may not be statistically significant. For example, several UEs in a cell could provide an abnormal cell profile if these UEs are all located at the cell edge. Thus, cell k is considered for diagnosis only if the condition $D_k > D_{th}$ is fulfilled, where D_k is the number of data sessions in cell k and D_{th} is a configurable threshold of the filter. Second, the cells are also filtered by their correlation value obtained in the previous step of the method. This filter would act as a detection stage in order to discard the cells whose profile is not similar to any of the RCPs and, as a consequence, they cannot be correctly diagnosed. It might happen that a cell profile is given by a high contribution of two or more fault classes. In such cases, the proposed system is not capable of providing a diagnosis and they should be analyzed by troubleshooting experts in order to determine whether the source of the problem is the same or not and, if not, to solve the problems independently. Thus, cell k is considered for diagnosis only if the condition $\rho_{k,ref} < \rho_{th}$ is satisfied, where ρ_{th} is a configurable threshold.

Lastly, once the cells have been filtered, the diagnosis at the cell level is computed. For cell k , the diagnosis is given by the RCP having the highest correlation to the cell profile of cell k . More specifically, the root cause c_i associated to that RCP is the diagnosis for the given cell.

5 Performance evaluation

In this section, the proposed method is evaluated and compared to other reference methods. To this end, dataset 1 is used due to its larger number of data sessions compared to the other datasets. In addition, the filters of the proposed system have the following configuration: $D_{th} = 8$ data sessions and $\rho_{th} = 0.70$.

The methods are evaluated by using a key metric in diagnosis referred to as Diagnosis Success Rate (DSR). It is estimated as the ratio between the number of well-diagnosed cells and the total number of cells. Note that, to calculate the proportion of observations correctly diagnosed, the real root cause, or “label” associated to the observed cell, must be known. This labeling has been carried out in this paper by applying the rules that troubleshooting experts use to diagnose faults in their networks. Such rules, manually derived from experience, are verified by means of the inspection personnel that is in charge of checking the status of the affected network elements, either remotely or on-site. Thus, these labeled data have been used as reference to calculate the diagnose success rate of the evaluated methods.

To benchmark the proposed method, the results are compared to baseline methods based on previous work. On the one hand, regarding the user diagnosis, the proposed method (i.e., a classification tree) is compared to other approaches that are based on the following machine learning techniques:

- The k -means clustering: as explained in the previous section, this technique is used to partition the data into a number of clusters. In this case, the data sessions (regardless of the serving cell) are the input samples to be divided into various clusters. Then, each cluster is analyzed by experts and labeled with the normal state or one of the considered faults as in the classification tree. In this way, users belonging to a certain cluster share the same diagnosis. The result of this process (i.e. the user diagnosis), which is equivalent to that obtained from applying the classification tree, can be used for cell diagnosis. The k -means clustering has previously been applied to fault detection in the context of wireless networks [46].
- The self-organizing map (SOM): it is a type of Neural Network used for clustering. The methodology is similar to the k -means algorithm, since both are example methods of unsupervised learning. The approach followed in this paper is explained in [8], with the difference that the diagnosis in [8] is carried out by using traditional KPIs, while in this case call traces are used.

On the other hand, regarding the cell diagnosis, the proposed method (that uses k -means to generate RCPs and Pearson's correlation coefficient to compute cell diagnosis) is compared to the method presented in [38]. Such a method requires manual adjustment of several thresholds and configurable parameters. Its operation is divided into two stages, as follows. In the first stage, problematic cells are detected by comparing the percentage of well-performing UEs to a certain threshold. Then, in the

second stage, the cells detected as faulty cells in the previous stage are diagnosed only if the percentage of UEs related to at least one fault is above a threshold. In that case, the selected fault is the one with the largest percentage of UEs having the same fault in the cell. Although the method is mainly based on selecting the most dominant fault (i.e., the one having the maximum percentage of UEs), note that the particular distribution of UEs (e.g., in Normal and Unknown categories) is key in the diagnosis process. Hereafter, the state-of-the-art method in [38] is referred to as ‘Max,’ as opposed to the correlation-based approach.

Note that the methods for user diagnosis can be interchangeably combined with the methods for cell diagnosis. All these combinations have been implemented for evaluation. Lastly, a reference method for cell diagnosis is based on applying the classification tree to traditional (cell level) KPIs. To do this, the average per cell of the selected metrics is firstly calculated. Then, the binary classification tree is applied to the obtained samples with the aim of determining the diagnosis. Note that the population Ω of the input space is composed of the set of cells included in the dataset rather than the carried data sessions. Thus, a difference with the proposed method is that only one fault class is associated to a given cell. On the contrary, with the proposed method, a distribution of data sessions over the fault classes is provided for each cell, namely the cell profile. A configuration parameter that both methods have in common is D_{th} in order to focus only on cells with a significant number of data sessions.

The first analysis is devoted to the user diagnosis. Figure 7 shows the DSR obtained in each category for the proposed (tree) and reference (k -means and SOM) methods. The ‘Normal’ category is not included in this figure since all data sessions with a throughput below th_1 are labeled (as real cause) with a faulty state, thus in the same way that the evaluated approaches. The last bars provide the global DSR, which is obtained by aggregating the UEs from each category (rather than averaging the DSRs). The values in brackets below each category represent the number of data sessions whose real diagnosis corresponds to such a category. According to these values, the dataset appears to be unbalanced, i.e. the number of samples varies greatly from one class to another. This is a consequence of using real data where some faults are more likely than others. It is observed that, from a perspective of the clustering techniques, the smallest class (“UL Int”) is more difficult to be identified than others. However, the accuracy of the evaluation should not be affected since the number of samples per class is in the order of hundreds. A closer inspection of Fig. 7 reveals that the proposed method provides good DSR for “DL Cong” since this problem is placed at the top of the tree. This means that, in situations where DL Cong

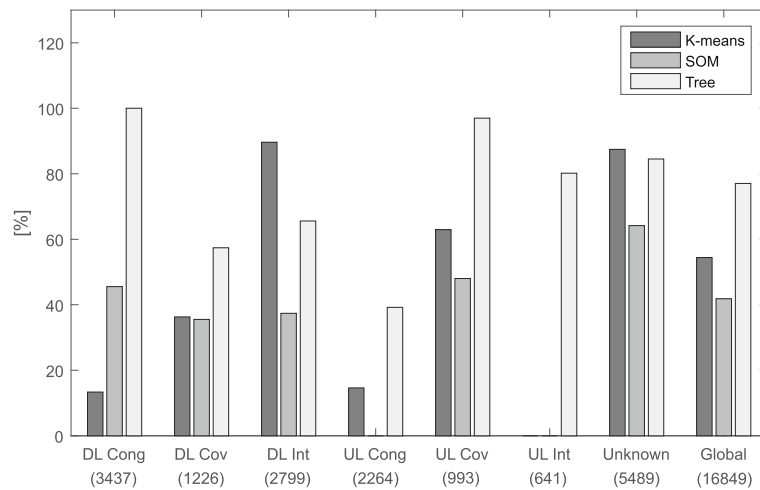


Fig. 7 Diagnosis success rate per class for user diagnosis methods

might be confused with other fault classes, the classification tree will select DL Cong as diagnosis. Another class with very high DSR for the proposed method is “UL Cov.” In this case, the good correlation between the metric and the associated fault has led to a higher DSR with respect to other classes, as occurred in the reference methods. It is also important that the methods provide high DSR for the Unknown class, since the addition of new faults in the diagnosis system will rely on the accuracy of this class. Lastly, regarding the last columns, it is noted that the proposed classification tree outperforms the reference methods in the global DSR. The main differences between the k -means and SOM methods are related to DL Cong and “DL Int” classes. Note that, since these two faults are closely related (i.e. a congestion problem usually entails higher interference levels), the bad DSR in DL Cong achieved by the k -means method is counteracted with a good DSR in DL Int. This issue is less pronounced in the case of the SOM method.

The existence of noise into the metrics may affect the performance of the user diagnosis. The PBD method mitigates the impact of the bias into the selected metrics in the same way as the purging methods. To analyze the impact of the noise variance, a synthetic noise has been added to the metrics. In particular, the noise is generated following the Additive white Gaussian noise (AWGN) model. The average of the AWGN is set to zero and the standard deviation is set to $0.2 \times \sigma$, where σ is the standard deviation of the metric data. Figure 8 shows the global DSR obtained for each method in presence/absence of AWGN. As observed, the greater the DSR is, the higher the impact of the noise is on the performance. Consequently, the most affected method is the proposed classification tree,

experiencing a decrease of about 5% with respect to the case without AWGN. However, the obtained value (71.6%) is still much greater than the value obtained by other methods.

With respect to intermittent perturbation, e.g., due to user mobility, the potential impact on the system performance is very limited due to temporal and spatial diversity factors. Specifically, in the time domain, the call traces are collected during a large period (typically, 15 min). This means that all the trace events and samples gathered during this period are considered for the elaboration of the metrics. Hence, this period is large enough to filter most propagation effects. In the spatial domain, all the users located in the service area of a cell contribute to determine its diagnosis. Thus, if one user experiences an eventual perturbation due to propagation effects, the users in other locations of the cell will counteract it.

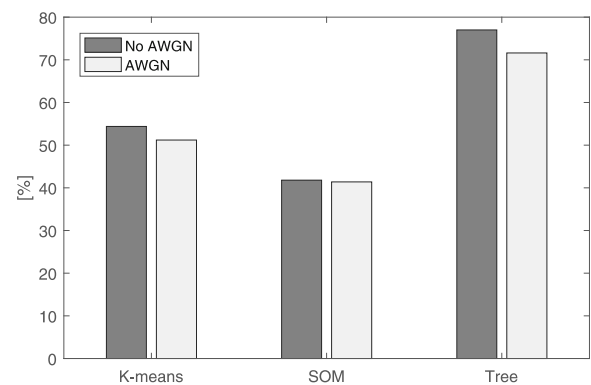


Fig. 8 Impact of noise on the global diagnosis success rate

The following analysis is related to the cell diagnosis, which is conducted based on the user diagnosis obtained in the previous evaluation. To do this, the methods “Max” and “Correlation” are combined with the previous methods, k -means, SOM and classification tree. Note that, in the case of solutions that comprise the proposed ‘Correlation’ method, a k -means clustering has been applied in order to find adequate RCPs for correlation. Figure 9 represents the DSR for such combined approaches. The method that applies the classification tree to cell-level KPIs is also included as a baseline. It is observed that the solutions based on reference clustering methods (i.e., k -means and SOM) provide worse performance than the classification tree with cell-level KPIs. This is mainly due to the bad user diagnosis that is obtained in the first stage of the solution. Moreover, the benefits of the proposed Correlation method are not evident using these solutions. If, on the contrary, the user diagnosis is carried out by means of the classification tree, the obtained DSR is better than applying the classification tree to cell-level KPIs. However, it is essential to use the proposed Correlation method in order to obtain a clear advantage over the reference method.

To further explain this, Fig. 10 shows the DSR obtained in each category for each method. As before, the values in brackets represent the actual number of data sessions in each category. It is observed that the amount of samples in some categories is scarce, so that the accuracy of the evaluation can be affected, especially for UL Int. The reason why Normal is not the dominant category is that cells with a low number of carried data sessions have been filtered. Due to the low number of served UEs, this group of cells are not critical from a troubleshooting perspective. Moreover, most of these cells belong to the

Normal category. Regarding the methods, it is noted that using the classification tree with cell-level KPIs leads to a worse diagnosis, especially for coverage and interference issues. One reason is that the performance of the classification tree is more sensible to threshold settings when applied to cell-level KPIs (instead of call traces), since only one sample per cell is used. Another reason is that, if the proportion of problematic UEs in the cell is small, the abnormal values are hidden when the average over all the UEs in the cells is calculated. This idea can also be applied to the ‘Max’ method, since in some cases the identification of faults is hindered by the dominance of ‘Normal’ class, having the maximum percentage of UEs. Therefore, it can be concluded that the proposed system, based on a classification tree and a correlation-based approach, provides more accurate diagnosis to troubleshooting experts than the baselines. In addition, under the SON framework, an improved DSR will increase the effectiveness of the subsequent stages in the process, such as the fault compensation.

To evaluate the capability of detection of the proposed method, the false positive and false negative rates have been calculated. In particular, Fig. 11 shows this information together with the number (in brackets) of data sessions that are classified as false positives or false negatives. As observed, the method based on a classification tree and a correlation-based approach provides a better trade-off between false positives and false negatives, because the number of cases is similar in both cases (i.e., 17 and 11). In addition, the number of false negatives is much lower than in other methods.

The gain in accuracy achieved by the proposed method is at the cost of increasing the operational complexity due to the management of a larger amount of information,

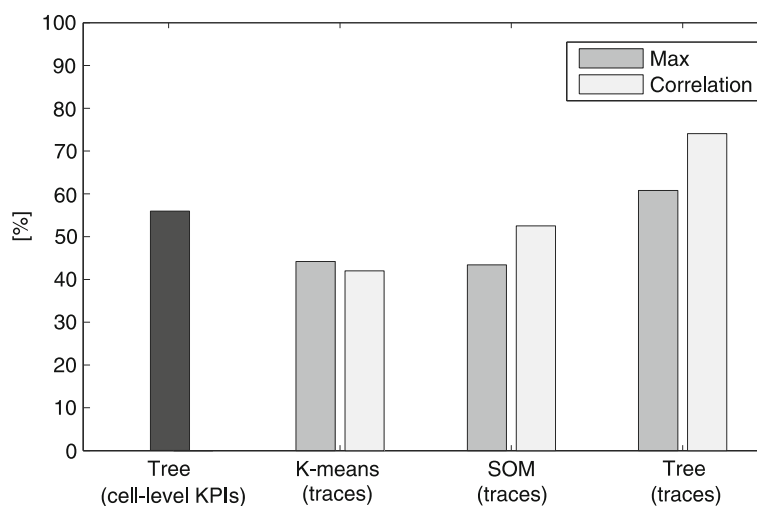


Fig. 9 Global diagnosis success rate for cell diagnosis methods

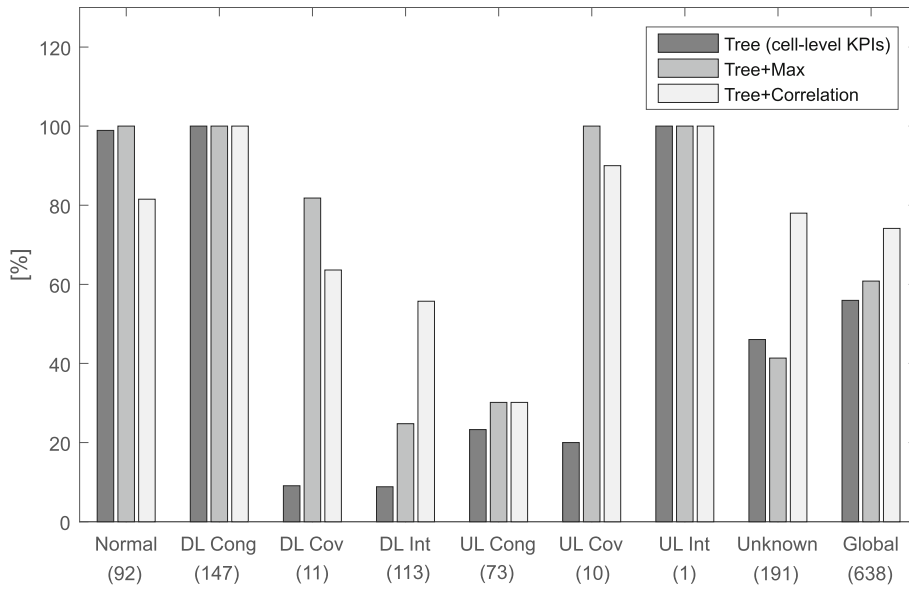


Fig. 10 Diagnosis success rate per class for cell diagnosis methods

given by the call traces. To estimate the computational cost, let C_{tree} be the cost of executing the classification tree. Its computation only requires elementary arithmetical operations, thus the complexity is $O(n)$. Let C_{profile} be the cost of generating the cell profile, whose complexity is also $O(n)$, since it mainly comprises data aggregation. The RCPs are not needed to be generated every time the cells are evaluated for diagnosis, e.g. they can be obtained from past evaluations in the same network. Thus, the computational cost of applying a clustering technique to find RCPs is not considered in this estimation. Let C_{corr} be the cost of calculating the Pearson's correlation coefficient for two data arrays. In this case, the computational cost is expected to be higher than C_{tree} due to the presence of the square root, whose complexity is $O(M(n))$ if the Newton's

method is used. Lastly, let N_{cell} be the number of cells in the dataset and N_{session} be the average number of data sessions per cell. Then, the complexity of the proposed method is estimated as $N_{\text{cell}} \cdot (N_{\text{session}} \times C_{\text{tree}} + C_{\text{profile}} + C_{\text{corr}})$, while in the case of using traditional KPIs (i.e. cell-level measurements), the estimated complexity is $N_{\text{cell}} \times C_{\text{tree}}$. Although there are substantial differences in the computational cost, however, with the current solutions of Big Data and Edge Computing, the required complexity is not an issue. In addition, the temporal restrictions associated to cell diagnosis are not the same as those applied to other network management functions such as packet scheduling at the link layer, which typically are more critical. In particular, the time to collect cell traces, which is defined by the Report Output Period (ROP), in general is much larger than the required time to execute the diagnosis method. The ROP is typically 15 min, while the proposed classification tree is executed in the order of milliseconds to seconds, depending on the number of cells. Thus, the rate at which the diagnosis system is fed with measurements is much slower than the required time to compute the diagnosis.

6 Conclusions

Mobile operators have focused their attention in improving the user satisfaction for data services in their networks. In an IP world, a user connection experiencing low throughput can be given by many factors. The RAN is commonly the critical segment of the end-to-end path that affects the service integrity due to the limited amount of resources available in the radio interface. The correlation between the user throughput and the related radio

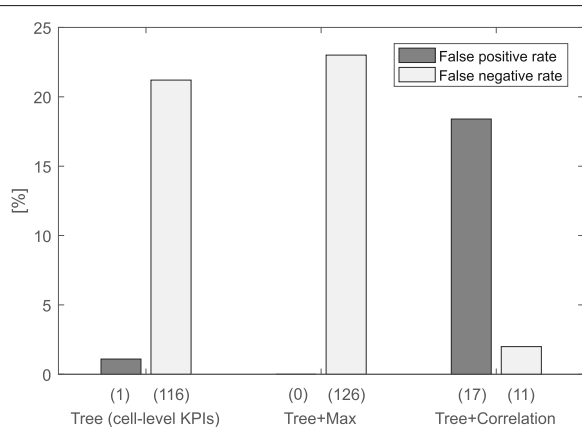


Fig. 11 False-positive and false-negative rates

conditions has been investigated in this paper with the aim of root-cause analysis. The study has been carried out through the use of large datasets collected from mature real networks.

More specifically, a diagnosis system has been developed to identify the most common radio causes that may impact user throughput. It takes advantage of the measurements that are collected per UE-basis as opposed to traditional counters and KPIs given on a cell-basis. The proposed system is based on a binary classification tree that is used to diagnose every data session in a cell. This tree allows the identification of data sessions suffering problems related to congestion, coverage or interference in either the downlink or the uplink from a set of metrics that have been carefully obtained considering some aspects such as the cell bandwidth and the data session duration. The diagnosis for each data session is then aggregated at the cell-level to provide a cell status (profile). After this, the set of cell profiles is then compared to several reference cell profiles in order to find anomalous cells. Such reference cell profiles has been previously generated by applying clustering techniques, e.g., the *k*-means. The comparison between cell profiles is made by calculating a correlation coefficient which is then used to determine the diagnosis of the cell.

Results show that the proposed system does not require a thorough adjustment of thresholds such as in other approaches, since the reference cell profiles automatically includes the particular characteristics of each network. Regarding the user diagnosis, the proposed method has been compared to other common clustering methods. It has been shown that the priority given by the proposed classification tree to the faults provides better accuracy in the diagnosis. With respect to the cell diagnosis, it has been shown that the performance of a small group of UEs experiencing poor radio conditions may be hidden by the dominant user performance in a cell. Due to this, the accuracy of the baselines is negatively affected. Conversely, the proposed method, by means of the generated RCPs, provides a better diagnosis for this kind of situations. Such an improvement leads troubleshooting experts to take appropriate recovery actions to solve the faults.

In addition, the proposed method can be applied to other radio access technologies since the thresholds for the metrics in the classification tree are calculated using percentiles. Considering that degradation is associated to extreme values of these metrics, the thresholds are automatically adapted to the radio access technology of the dataset. However, beyond the particular range of the metrics, it is important that the metrics are equivalent to those utilized in this paper in order to detect the considered faults. For example, to identify the problem of lack of coverage, a metric related to the signal strength should be used.

Further work is required to identify, in an automated manner, what values of throughput are considered to be abnormally low in a certain network based on other metrics. In particular, to calculate the threshold for the throughput metric, automatic learning techniques can be applied to the datasets, replacing the sensitivity analysis.

Acknowledgements

This work has been partially funded by Optimi-Ericsson, Junta de Andalucía (Agencia IDEA, Consejería de Ciencia, Innovación y Empresa, ref. 59288, and Proyecto de Investigación de Excelencia P12-TIC-2905) and ERDF.

Authors' contributions

All authors contributed extensively to the work presented in this paper. PM designed the proposed algorithm, performed the experiments, wrote the paper. RB supervised the design of the algorithm and the performance evaluation, wrote the paper. EC was in charge of the data collection and development of the supporting tools to process the data; he revised the paper. AG designed the reference algorithms, analyzed the results, wrote the paper. EJK designed the reference algorithms, analyzed the results, wrote the paper. NF formulated the problem, supervised the design of the algorithm, revised the paper. All authors read and approved the final manuscript.

Competing interests

The authors declare that they have no competing interests.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Author details

¹Communications Engineering Dept., University of Málaga, Málaga, Spain.

²Ericsson, 29590 Campanillas, Málaga, Spain.

Received: 29 April 2016 Accepted: 6 July 2017

Published online: 20 July 2017

References

1. Cisco, Visual networking index: global mobile data traffic forecast update. Tech. Rep. Cisco, 2014–2019 (2015)
2. 3GPP, Self-organizing networks (SON): concepts and requirements, version 12.1.0, TS 32.500 (2014–12)
3. 3GPP, Evolved Universal Terrestrial Radio Access (E-UTRA); Physical layer; Measurements, version 12.2.0, TS 36.214 (2015–03)
4. 3GPP, Key Performance Indicators (KPI) for Evolved Universal Terrestrial Radio Access Network (E-UTRAN): definitions, version 12.0.0, TS 32.450 (2014–10)
5. J Laiho, K Raivio, P Lehtimäki, K Hatanen, O Simula, Advanced analysis methods for 3G cellular networks. *IEEE Trans. Wirel. Commun.* **4**(3), 930–942 (2005)
6. Y Zhang, A Årvidsson, in *Proc. of 2012 ACM SIGCOMM Workshop on Cellular Networks: Operations, Challenges, and Future Design*. Understanding the characteristics of cellular data traffic (ACM, Helsinki, 2012), pp. 13–18
7. EJ Khatib, R Barco, I Serrano, P Muñoz, in *Proc. of IEEE GLOBECOM 2014*. LTE performance data reduction for knowledge acquisition (IEEE, Austin, 2014), pp. 270–274
8. A Gómez-Andrades, P Muñoz, I Serrano, R Barco, Automatic root cause analysis for LTE networks based on unsupervised techniques. *IEEE Trans. Veh. Technol.* **65**(4), 2369–2386 (2015)
9. 3GPP, Telecommunication management; Subscriber and equipment trace; Trace concepts and requirements, version 13.0.0, TS 32.421. (2015-06)
10. 3GPP, Universal terrestrial radio access (UTRA) and evolved universal terrestrial radio access (E-UTRA); Radio measurement collection for minimization of drive tests (MDT); Overall description; Stage 2, version 12.2.0, TS 37.320 (2014-09)
11. LC Schmelz, et al., in *Wireless World Research Forum Meeting 20*. Self-configuration, -optimisation and -healing in wireless networks (WWRF, Stockholm, 2008)

12. R Barco, P Lázaro, V Wille, L Díez, S Patel, Knowledge acquisition for diagnosis model in wireless networks. *Expert Syst. Appl.* **36**(3), 4745–4752 (2009)
13. R Barco, P Lázaro, P Muñoz, A unified framework for self-healing in wireless networks. *IEEE Commun. Mag.* **50**(12), 134–142 (2012)
14. GF Ciocarlie, CC Cheng, C Connolly, U Lindqvist, S Nováczki, H Sanneck, M Naseer-ul-Islam, in *Proc. of 11th International Symposium on Wireless Communications Systems (ISWCS)*. Managing scope changes for cellular network-level anomaly detection (IEEE, Barcelona, 2014), pp. 375–379
15. J Johansson, WA Hapsari, S Kelley, G Bodog, Minimization of drive tests in 3GPP release 11. *IEEE Commun. Mag.* **50**(11), 36–43 (2012)
16. J Puttonen, J Turkka, O Alanen, J Kurjenniemi, in *Proc. of IEEE 21st International Symposium on Personal Indoor and Mobile Radio Communications (PIMRC)*. Coverage optimization for minimization of drive tests in LTE with extended RLF reporting (IEEE, Istanbul, 2010), pp. 1764–1768
17. A Galindo-Serrano, B Sayrac, S Ben Jemaa, J Riihijarvi, P Mahonen, in *Proc. of 8th International Conference on Cognitive Radio Oriented Wireless Networks (CROWNCOM)*. Harvesting MDT data: Radio environment maps for coverage analysis in cellular networks (IEEE, Washington, 2013), pp. 37–42
18. RU Mondal, T Ristaniemi, J Turkka, in *Proc. of International Conference on Localization and GNSS (ICL-GNSS)*. Genetic algorithm optimized grid-based RF fingerprint positioning in heterogeneous small cell networks (IEEE, Gothenburg, 2015)
19. PC Lin, in *Proc. of IEEE 3rd Global Conference on Consumer Electronics (GCCE)*. Minimization of drive tests using measurement reports from user equipment (IEEE, Tokyo, 2014), pp. 84–85
20. F Chernogorov, J Puttonen, in *Proc. of IEEE 24th International Symposium on Personal Indoor and Mobile Radio Communications (PIMRC)*. User satisfaction classification for Minimization of Drive Tests QoS verification (IEEE, London, 2013), pp. 2165–2169
21. F Chernogorov, T Nihtilä, in *Proc. of IEEE 75th Vehicular Technology Conference (VTC Spring)*. QoS verification for minimization of drive tests in LTE networks (IEEE, Yokohama, 2012)
22. MP Wylie-Green, T Svensson, in *Proc. of IEEE Global Telecommunications Conference (GLOBECOM 2010)*. Throughput, capacity, handover and latency performance in a 3GPP LTE FDD field trial (IEEE, Miami, 2010)
23. B McWilliams, Y Le Pezennec, G Collins, in *Proc. of XVth International Telecommunications Network Strategy and Planning Symposium (NETWORKS)*. HSPA+ (2100 MHz) vs LTE (2600 MHz) spectral efficiency and latency comparison (IEEE, Rome, 2012)
24. L Zhang, T Okamawari, T Fujii, in *Proc. of IEEE 75th Vehicular Technology Conference (VTC Spring)*. Performance evaluation of end-to-end communication quality of LTE (IEEE, Yokohama, 2012)
25. V Buenestado, JM Ruiz-Avilés, M Toril, S Luna-Ramírez, A Mendo, Analysis of throughput performance statistics for benchmarking LTE networks. *IEEE Commun. Lett.* **18**(9), 1607–1610 (2014)
26. A Imran, A Zoha, Challenges in 5G: how to empower SON with big data for enabling 5G. *IEEE Netw.* **28**(6), 27–33 (2014)
27. N Baldo, L Giupponi, J Mangués-Bafalluy, in *Proc. of 20th European Wireless Conference*. Big data empowered self organized networks (VDE, Barcelona, 2014)
28. EJ Khatib, R Barco, P Muñoz, I de la Bandera, I Serrano, Self-healing in mobile networks with big data. *IEEE Commun. Mag.* **54**(1), 114–120 (2016)
29. P Szilagy, S Novaczki, An automatic detection and diagnosis framework for mobile communication systems. *IEEE Trans. Netw. Serv. Manag.* **9**(2), 184–197 (2012)
30. R Barco, V Wille, L Díez, M Toril, Learning of model parameters for fault diagnosis in wireless networks. *Wirel. Netw.* **16**(1), 255–271 (2010)
31. L Bennacer, L Ciavaglia, A Chibani, Y Amirat, A Mellouk, in *Proc. of IEEE Network Operations and Management Symposium (NOMS)*. Optimization of fault diagnosis based on the combination of Bayesian Networks and Case-Based Reasoning (IEEE, Maui, 2012), pp. 619–622
32. EJ Khatib, R Barco, A Gómez-Andrades, I Serrano, Diagnosis based on genetic fuzzy algorithms for LTE self-healing. *IEEE Trans. Veh. Technol.* **65**(3), 1639–1651 (2016)
33. Y Zhang, N Liu, Z Pan, T Deng, X You, in *Proc. of IEEE/CIC International Conference on Communications in China (ICCC)*. A fault detection model for mobile communication systems based on linear prediction (IEEE, Shanghai, 2014), pp. 703–708
34. P Muñoz, R Barco, I Serrano, A Gómez-Andrades, Correlation-based time-series analysis for cell degradation detection in SON. *IEEE Commun. Lett.* **20**(2), 396–399 (2016)
35. B Cheung, G Kumar, S Rao, Statistical algorithms in fault detection and prediction: toward a healthier network. *Bell Labs Tech. J.* **9**(4), 171–185 (2005)
36. N Samaan, A Karmouch, Network anomaly diagnosis via statistical analysis and evidential reasoning. *IEEE Transactions on Network and Service Management.* **5**(2), 65–77 (2008)
37. A Gómez-Andrades, R Barco, I Serrano, P Delgado, P Caro-Oliver, P Muñoz, Automatic root cause analysis based on traces for LTE self-organising networks. *IEEE Wirel. Commun. Mag.* **23**(3), 20–28 (2016)
38. A Gómez-Andrades, R Barco, P Muñoz, I Serrano, Data analytics for diagnosing the RF condition in Self-Organising Networks. *IEEE Trans. Mob. Comput.* In press, (2016)
39. J Turkka, F Chernogorov, K Brigatti, T Ristaniemi, J Lempiäinen, An approach for network outage detection from drive-testing databases. *J. Comput. Netw. Commun.* **2012**, 1–13 (2012)
40. A Zoha, A Saeed, A Imran, M Imran, A Abu-Dayya, in *Proc. of 11th International Conference on the Design of Reliable Communication Networks (DRCN)*. Data-driven analytics for automated cell outage detection in Self-Organizing Networks (IEEE, Kansas City, 2015), pp. 203–210
41. O Onireti, A Zoha, J Moysen, A Imran, L Giupponi, M Imran, A Abu Dayya, A cell outage management framework for dense heterogeneous networks. *IEEE Trans. Veh. Technol.* **65**(4), 2097–2113 (2015)
42. H Holma, A Toskala, *LTE for UMTS: Evolution to LTE-Advanced*. (John Wiley & Sons Ltd, The Atrium, Southern Gate, Chichester, 2011)
43. K Aho, O Alanen, J Kaikkonen, in *Proc. of 10th International Conference on Networks (ICN)*. CQI reporting imperfections and their consequences in LTE networks (IARIA, St. Maarten, 2011), pp. 241–245
44. 3GPP, Evolved Universal Terrestrial Radio Access (E-UTRA); Layer 2 - Measurements, version 12.0.0, TS 36.314 (2014–09)
45. R Khanafer, B Solana, J Triola, R Barco, L Moltsen, Z Altman, P Lazaro, Automated diagnosis for UMTS networks using Bayesian network approach. *IEEE Trans. Veh. Technol.* **57**(4), 2451–2461 (2008)
46. M Wazid, A Das, An efficient hybrid anomaly detection scheme using K-means clustering for wireless sensor networks. *Wirel. Pers. Commun.* **90**(4), 1971–2000 (2016)

Submit your manuscript to a SpringerOpen[®] journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at ► springeropen.com