

Anomaly Detection and Modeling in 802.11 Wireless Networks

Anisa Allahdadi · Ricardo Morla

the date of receipt and acceptance should be inserted later

Abstract IEEE 802.11 Wireless Networks are getting more and more popular at university campuses, enterprises, shopping centers, airports and in so many other public places, providing Internet access to a large crowd openly and quickly. The wireless users are also getting more dependent on WiFi technology and therefore demanding more reliability and higher performance for this vital technology. However, due to unstable radio conditions, faulty equipment, and dynamic user behavior among other reasons, there are always unpredictable performance problems in a wireless covered area. Detection and prediction of such problems is of great significance to network managers if they are to alleviate the connectivity issues of the mobile users and provide a higher quality wireless service. This paper aims to improve the management of the 802.11 wireless networks by characterizing and modeling wireless usage patterns in a set of anomalous scenarios that can occur in such networks. We apply time-invariant (Gaussian Mixture Models) and time-variant (Hidden Markov Models) modeling approaches to a dataset generated from a large production network and describe how we use these models for anomaly detection. We then generate several common anomalies on a Testbed network and evaluate the proposed anomaly detection methodologies in a controlled environment. The experimental results of the Testbed show that HMM outperforms GMM and yields a higher anomaly detection ratio and a lower false alarm rate.

Keywords 802.11 Access Points · Network Usage · Gaussian Mixture Models · Hidden Markov Models · Anomaly Detection

1 Introduction

Wireless 802.11 networks are getting more and more popular in providing Internet access for a large number of users in university campuses, enterprises, urban areas, and many other public places. These large-scale networks, particularly speaking of IEEE 802.11 Infrastructure mode, consist of basic network components: Wireless Stations, wired stations, and the Access Points (AP) that function as connection links between the wired and wireless sections. The APs provide coverage and capacity for supporting mobile clients with heterogeneous devices and a variety of applications. Among the many characteristics of such large-scale network is the transition of huge volumes of traffic as a result of intensive usage from different locations all over the coverage area. The mobile clients demand reliable connection and high performance in all circumstances and expect their applications to work smoothly around the wireless covered field, but this is an ideal case which is not always achievable. The wireless users, most of the time suffer from low coverage, intermittent connectivity, authentication failure, degraded performance and many other complications originated from the unreliable nature of wireless connection and dynamic usage pattern of other users in the vicinity.

Having further explored the connectivity procedure in wireless networks, some inherent concerns and dilemmas become more clear. In Wireless 802.11 networks, mobile stations perform an active or passive scanning process to discover available APs in the vicinity and

A. Allahdadi
INESC TEC, Faculty of Engineering, University of Porto,
Campus da FEUP, Rua Dr. Roberto Frias N8 378,
4200-465 Porto, Portugal
E-mail: anisa.allahdadi@inescporto.pt

R. Morla
E-mail: ricardo.morla@fe.up.pt

connect to an AP with the highest received signal strength (RSS) [24]. This association strategy, only based on RSS, can lead to many connectivity problems and performance issues as it may result in significant load imbalance between APs. The overloaded APs can still present high RSS and try to accommodate more stations while other APs are only slightly loaded or even idle. Another source of performance degradation in WLANs is the multi-rate flexibility and the fairness mechanism of the MAC protocol- when a station far from the AP reduces its bit rate to avoid repeated unsuccessful frame transmission and as a result degrades the throughput of the other stations associated with the same AP [17]. In addition to the aforementioned problems, due to the unreliable and time-varying nature of the wireless channels, 802.11 networks usually suffer from many pitfalls such as exposed and hidden terminals, capture effect, interferences, signal fading, inconsistent coverage, and many other examples. In such circumstances, high packet loss is observed [13]- that results in inconsistent connectivity and low performance. Network managers are concerned about discovering such sort of problems and abnormal events that occur in their network. Detection of anomalies is not only advantageous for prompting immediate administrative actions but also useful for long-term network design, planning, and maintenance decisions as the network infrastructure and usage evolve over time.

In large deployments of 802.11 networks with varying usage, channel conditions, and operational constraints, network managers often demand tools that provide them with a comprehensive view of the entire network for automatic detection of the problems. In such widespread networks, where at any moment there is a high possibility of mal-functioning of APs and user devices, the necessity of such automatic tools or applications is vital to preserve the quality of service at an acceptable level. Monitoring the infrastructure by any means rather than intelligent diagnostic tools seems inconvenient in practice or overpriced in budget. For example, it is expensive to deploy third party devices like sensors and sniffers individually on clients machines or APs for detection of problems in different OSI layers, as studied earlier [7, 12, 25]. And it seems impractical for network staff to walk around the wireless covered area with a device in their hand monitoring the network and measuring the quality of connections at any time. In this paper, we propose an automatic diagnostic tool that analyzes the usage data of the APs- collected from a RADIUS authentication server. We apply probabilistic learning algorithms to produce a model for each access point or group of access points, and identify anomalous events with a margin of certitude. AP usage modeling and

anomaly detection in hotspots would assist network administrators to ensure long-term quality of service by analyzing various connectivity factors of wireless users in particular localities. We propose probabilistic graphical model- and in particular HMM- to establish a comprehensive image of the evolving structure of wireless networks, to distinguish usage behaviors in different locations and grouping context and their correlations and dependencies, and to represent the spatio-temporal anomalous patterns detected in wireless networks. In the current work we focus more on proposing individual models for APs as the ground truth data is only available through the single AP Testbed deployment and the multiple APs' experiment is planned for the future work. The prospective methodology is based on the development of HMM models and a detection tool using WiFi campus data; our recent contributions [9, 10] have taken this approach into account. As a preliminary investigation on the subject, we focused on short 802.11 sessions recorded through RADIUS authentication as a network artifact and an indicator of quality of wireless access [9]. In [10] an exhaustive analysis is performed for outlier detection in 802.11 wireless networks using HMM variations- single HMM, mixture of HMMs and individual HMMs- and is evaluated by the state of the art statistical methodologies. Furthermore a number of network anomalous patterns are represented, in the same study, considering HMM parameters such as hidden states' transition and partial likelihood of the observation sequences. In the present study we considered HMM and its counterpart time-invariant methodology- Gaussian Mixture Model (GMM)- to investigate the temporal relevancy of the employed data, whether a simpler time-invariant model such GMM is adequate to detect anomalies or a more complex model like HMM is really needed. These two methodologies are analyzed and compared with each other both in modeling and anomaly detection experiments.

This paper contains two main parts: 1) analysis and modeling of 802.11 AP usage and exploring the time dependency of the employed data, and 2) identification and detection of different types of anomalies and characterizing them efficiently. The aforementioned objectives are investigated on a large dataset of AP usage and examined on a smaller scale testbed for the purpose of evaluation.

The rest of the paper proceeds as follows. In section 2, the related work and the most recent researches relevant to the current work are presented. In section 3, the wireless setup procedure in infrastructure mode is characterized and the key attributes and functionalities of RADIUS protocol are defined. Section 4 deals with the process of data accumulation as a result of wireless

users' association attempts, and presents a set of main features extracted from the dataset and feature selection techniques for further analysis. Statistical modeling of the AP usage data, categorized as time-invariant and time-variant approaches are provided in section 5 and a brief discussion is enclosed at the end of the section comparing these methodologies. In section 6, we describe how the proposed models serve to detect and characterize different forms of anomalies. In section 7, the experimental results are analyzed and discussed, and the evaluation process is represented based on the deployed testbed in a controlled environment. In section 8, the major conclusions are provided and prominent directions for future work are identified.

2 Related Work

2.1 Wireless Measurement Tools

Several prior works are dedicated to studying the dynamics of wireless network behavior, as well as the performance and reliability of WLAN technologies [12] [31] [30] [20]. In [12] a system called Jigsaw is presented which uses multiple monitors to provide a single unified view of physical, link, network and transport-layer activities, including inference techniques for the particular issues of 802.11. The authors deployed an infrastructure with over 150 radio monitors that capture 802.11b and 802.11g activities in a university building to investigate the causes of performance degradation. Significant challenges of such vast distributed monitoring system include the necessity of hardware and software instrumentations on each and every monitor and the scalable synchronization difficulties and inaccuracies. For this reason, most wireless management techniques avoid broad modifications in the clients devices, sensors, sniffers and monitors deployed in the large wireless covered area.

In another line of research a Passive Interference Estimator (PIE) is presented [31] which provides a fine-grained estimation of link interferences in WLAN. PIE provides an estimate of WLAN interference caused by client mobility, dynamic traffic loads, and varying channel conditions. This work is inspired by two previous WLAN monitoring approaches: the aforementioned Jigsaw [12] and WIT [21]. The PIE producers use sniffing at APs to avoid deploying additional monitors similar to Jigsaw, but with the penalty of missing a portion of uplink client traffics and hence uplink client conflicts. However, they proposed an accurate approach in estimating link interference by providing a conflict graph in real time.

In a similar direction of work, fine-grained detection algorithms are proposed that are able to distinguish the root-causes of performance degradation at the physical layer [30]. It is described that various faults, such as hidden terminals, capture effects and noise, could have the same propagation effects on the network layer (degraded throughput) and therefore could lead to the same remediation techniques from 802.11 (rate fallback), while they have completely different origins in the physical layer. Hence, the researchers of this work designed a unified framework for this purpose, called MOJO, that combines the observations from multiple distributed sniffers and diagnoses the granularity of the root causes to suggest appropriate remedies for different physical faults. Although the proposed framework measures the impact of the most commonly observed faults on different network layers, it is still a client side monitoring system and suffers from the extensive sniffer distribution all over the wireless covered area.

In [20], WiMed is proposed that uses only local measurements from commodity 802.11 NICs for understanding how the medium is utilized, and for inspecting the causes of interferences (including non-802.11 devices). WiMed provides a time-domain view of how the medium is used in a given 802.11 channel, and identify the root causes of interference using physical layer properties such as bit error patterns and medium busy times. The authors refrain from elaborate instrumentation and dedicated infrastructure, however detectors are only implemented for interference and contention, and there is a higher confidence for recognition of non-802.11 interferer rather than 802.11 sources of interference.

All the above literatures expose the difficulties in monitoring the wireless environment thoroughly, and the challenges of performance estimation in such complex networks. Most cases- require heavy instrumentation of the user devices and focus on specific anomalies affecting individual users- thus neither considering usage trend nor location related anomalies.

2.2 Usage Modeling and Anomaly Detection

There are several lines of research that take an approach closely related to our work. In [23] AP usage and daily keep-alive events of mobile stations in 802.11 hotspots in infrastructure mode are analyzed and modeled. In this work, generative probabilistic models are investigated such as Gamma mixture of exponentials and Conditional probability models considering dependencies between consecutive samples in time. The generative statistical models and experimental results of this work conducted on a very similar dataset to ours

provide some broad insight into AP usage and illustrate those aspects of such networks that benefit our work.

In [22], a usage pattern called "abrupt ending" is explored in a similar dataset, and it concerns the disassociation of a large number of wireless sessions in the same AP within a one second window, or in a nutshell "simultaneous session ending". The authors of this work, further investigate this concept and introduce some anomalous patterns that might be in correlation with the occurrence of this phenomena. For instance, they propose that interference across the AP vicinity could be deduced when abrupt endings happen to neighboring APs within specified time interval, or the AP overloaded could be inferred when the continued sessions are present after abrupt endings. There are a number of other anomalous patterns reported in this paper such as AP halt/crash, persistence interference and intermittent connectivity. The classification and analysis of these anomaly-related patterns performed in this research, inspired our work to regenerate similar anomalies in a real Testbed to experiment and evaluate the HMM methodologies practiced in the current study.

2.3 HMM Applications in Network Analysis

In wireless networking, HMMs are employed to address various aspects of network measurement and analysis. Hierarchical and Hidden Markov based techniques are analyzed in [19] to model 802.11b MAC-to-MAC channel behavior in terms of bit error and packet loss. The authors employed two random variables in packet loss process, inter-arrival-rate and burst-length of packet loss, and applied the traditional two-state Markov chain. The results demonstrates that two-state Markov chain provides an adequate model for the 802.11b MAC-to-MAC packet loss process. Furthermore, in regard to bit error modeling, three other Markov-based chains are evaluated: full-state, hidden, and hierarchical Markov chains. Among these chains, it is illustrated that the full-state Markov bit error model of order 9 and above, renders the best performance. Since the main concern to use HMM in this example is to generate error traces, a simple three-state HMM is designed and utilized for one HMM solution: the adjustment of model parameters to best account for the observed signal.

In a more recent line of research in [18] a multi-level approach involving HMMs and Mixtures of Multivariate Bernoullis (MMB) is proposed to model the long and short time scale behavior of wireless sensor network links, that is, the binary sequence or trace of packet receptions (1s) and losses (0s) in the link. In this approach, HMM is applied to model the long-term evolution of the trace, and the hidden states correspond to

packet reception rate. Within the aforementioned hidden states, the short-term evolution of the trace is modeled by either another HMM or by a MMB. That is how the multilevel, or in this case the two level approach, is formed. The notion of multilevel HMM, or higher dimensional HMM, is an impressive concept regarding to our own work, and we intend to make use of this approach to improve our HMM variations for anomalous pattern detection in the future work.

One of the salient applications of HMMs addressed in wireless networking is prediction. For instance in [11] HMMs are utilized to model and predict the spectrum occupancy of sharing radio bands. The channel status prediction is considered as a binary series prediction problem, as channel occupancy can be represented as idle or busy depending on the presence or absence of a primary user activity. An ergodic two-state discrete HMM deals with this problem. Some other prominent work has been done on a very similar subject in radio spectrum sensing and status prediction using HMMs in [8, 16, 32].

Furthermore, in another related work, HMMs are applied for modeling and prediction of user movement in wireless networks to address issues in Quality of Service (QoS) [26]. User movement from an AP to an adjacent AP is modeled using a second-order HMM. Although the authors demonstrated the necessity of using HMM instead of Markov chain model, the proposed model is only practical for small wireless networks with a few number of APs, not huge enterprises or widespread campuses.

As the above literatures indicate and to the best of our knowledge, HMM related studies in wireless network management are rarely used specifically in performance anomaly detection.

3 Wireless Setup in Infrastructure Mode

In this section we describe how a 802.11 station associates to an access point and how our setup authenticates the user and authorizes access to the network.

3.1 Association of Wireless Station to Access Point

The process of the association of a wireless mobile station to an AP, as it is currently implemented by most manufacturers is described as follows: A wireless station scans the available channels of each AP in the neighborhood and listens to the beacon (passive approach) or probe response frames (active approach). IEEE 802.11 protocol defines a number of Wi-Fi channels ranging from 2.4 GHz to 5.9 GHz. The Wi-Fi channels that are

Table 1 The Key Attributes of RADIUS Accounting Table

Acct-Status-Type	has three values: Start, Alive and Stop. A Start record is created when a user session begins. An Alive record is registered after each 10 or 15 minutes for the users that are still connected. A Stop record is generated when the session ends.
Acct-Session-Id	is a unique number assigned to each session to facilitate matching the Start and Stop records in a detail file, and to eliminate duplicate records.
Acct-Session-Time	records the user's connection time in seconds. This information could be included in Alive or Stop records.
Acct-Delay-Time	is the number of seconds passed between the event and the current attempt to send the record. The approximate time of an event can be determined by subtracting the Acct-Delay-Time from the time of the record's arrival on the RADIUS accounting server.
Called-Station-Id & Calling-Station-Id	record the IP address of the AP (Called Station) and the wireless user (Calling Station) connected to that AP.
Timestamp	records the time of arrival on the RADIUS Accounting host measured in seconds since the epoch (00:00 January 1, 1970). It provides a machine-friendly version of the logging time at the beginning of the accounting record.
Acct-Input-Octets & Acct-Output-Octets	records the number of bytes received (Acct-Input-Octets) and sent (Acct-Output-Octets) during a session. These values appear in Alive or Stop records.
Acct-Input-Packets & Acct-Output-Packets	records the number of packets received (Acct-Input-Packets) and sent (Acct-Output-Packets) during a session. These values appear in Alive or Stop records.

the concern of this work (802.11 b/g/n) are listed in the 2.4 GHz range and consist of one to eleven channels (up to fourteen in some countries). The wireless station stores the received signal strength indicator (RSSI) of the APs in the vicinity and other relevant information such as extended service set identification (ESSID), encryption type (e.g. WPA, WEP), etc. When the scanning process is over, the wireless station selects an AP with the highest RSSI among the observed APs in its proximity. After the process of authentication/ authorization is accomplished, the permission is granted to the wireless station and the connection is established. Forthwith, the wireless station is associated with the new AP and the user is ready to send and receive traffic through that AP. The wireless station will be disassociated from the current AP under the mobility circumstances, AP shutdown or halt, RSSI recession or some other normal or abnormal consequences of network fluctuations. The process of AP selection only based on the strongest RSSID lead to aforesaid load imbalance problem, while some APs are overcrowded and the other available APs remain idle.

3.2 Remote Authentication Dial-In User Service (RADIUS)

Remote Authentication Dial-In User Service (RADIUS) is a network protocol that enables remote access servers to communicate with a central server to authenticate dial-in users and authorize their access to the requested system or service. RADIUS is commonly used by Internet Service Providers (ISPs), cellular network providers, and corporate and educational networks, and it allows the management of user profiles in a central database

that all remote servers can share. Having a central service facilitates the process of tracking usage for billing and network statistics. RADIUS is a de facto industry standard used by a number of network product companies and it is a proposed IETF standard. This protocol is used to provide network authentication, authorization, and accounting services, and it is particularly described in Request for Comments (RFC) 2865 and RFC 2866.

According to RADIUS protocol, whenever a client associates to an 802.11 AP, a log event "START" is recorded in the accounting database. While the client is still connected to this AP, every 10 or 15 minutes (based on the server configuration) an interim log event "ALIVE" is issued to refresh the connection between the client and the AP. Eventually, when the user decides to disconnect from the network, or for some reason it is forced to leave the network, a log event "STOP" is recorded, which marks the end of the association period of this user. Each log record includes some key attributes of time-stamp, session ID, association duration, number of input and output packets/octetets. Table 1 present a brief explanation of some of these key attributes more relevant to this work.

RADIUS serves three main functionalities:

- Authenticates users before granting them access to the network.
- Authorizes the authenticated users for specific network services.
- Accounts the usage activity of the authorized users for the services in use.

AAA stands for "Authentication, Authorization, and Accounting". It defines an architecture that authenticates and grants authorization to users and accounts for

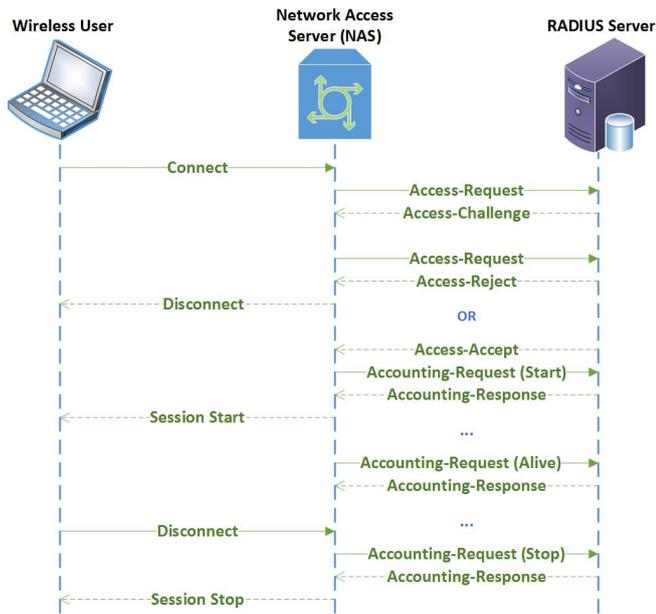


Fig. 1 The Authentication and Authorization Process in RADIUS

their activity. When AAA is not used, the architecture is described as "open", where anyone can gain access and do anything, without any tracking.

3.2.1 Authentication and Authorization

Authentication refers to the process of validating the identity of the user by matching the credentials provided by the user on the AAA server. If the credentials match, the user is authenticated and gains access to the network. On the contrary, if the credentials mismatch, authentication fails and network access is denied. Authentication can also fail, due to user incorrectly entering the credentials. A network administrator can choose to permit limited network access to unknown users, for instance the guests of a conference or a temporary public event in academic environments.

Authorization deals with the process of deciding what permissions are granted to the user. For example, the user may or may not be permitted certain kinds of network access or allowed to issue certain commands. Typically, a user login consists of a query (Access-Request) from the NAS to the RADIUS server and the RADIUS server either grants or denies authorization (Access-Accept or Access-Reject) based on the information passed by in the request query. In each case, the RADIUS server manages the authorization policy and the NAS enforces the policy. The process of authentication and authorization is delineated in Figure 1.

3.2.2 Accounting

Accounting refers to the recording of resources users consume during the time they are connected to the network. The information gathered can include the total system time used, and the amount of data sent or received by the user during a session. Over a network session, the NAS periodically sends an accounting data of user activity to the server (in "Alive" or "Stop" sessions). This data is mainly used for the billing purposes. However, we used the accounting information for the reason of network monitoring and management as the log dataset is already stored in a central database, the RADIUS server, and facilitates the data collection process.

The detailed information of users' activities is not included in the summary sent by NAS- for instance the visited web sites or particular protocols in use is local to the NAS- and is not available to the RADIUS server. Transactions between the client and RADIUS server are authenticated through the use of a shared secret, which is never sent over the network. In addition, user passwords are sent encrypted between the client and RADIUS server to eliminate the possibility of snooping on an insecure network.

4 Data Description and Feature Selection

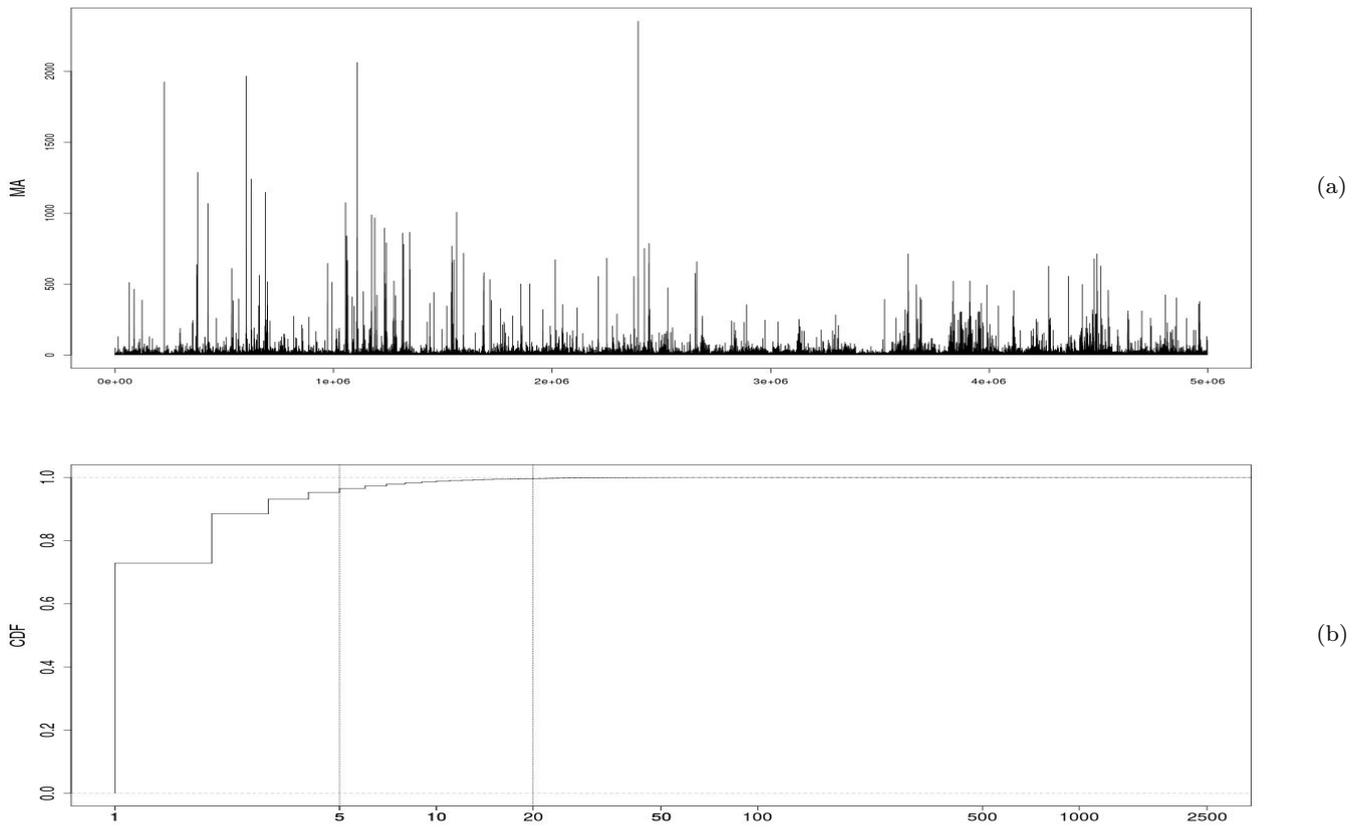
In this section we present the main dataset used in this paper, provide some preliminary statistical analysis and describe the key features emerge from the raw dataset as well as the of process of feature selection for modeling and further investigations.

4.1 Large Dataset

For the current study, we use RADIUS authentication log data collected at the hotspot of the Faculty of Engineering of the University of Porto (FEUP). The University hotspots are part of the Eduroam European wireless academic network initiative. The trace data consists of the daily summary of connections between 364 APs and their corresponding wireless stations collected in almost two years, from January 1, 2010 to December 22, 2011. The university campus contains over 30 buildings, including classrooms, administrative offices, auditoriums, libraries, cafeterias, laboratories, etc. During the mentioned period, the usage record of more than 45 thousand users was observed through the established connections of over 24 million sessions. Table 2 depicts

Table 2 The semester-level evolution of hotspot usage during two years

Academic Semesters	# APs	# Users	# Sessions	Total Input Traffic (TB)	Total Output Traffic (TB)
Spring 10/11	238	15564	5127823	148	253
Fall 10/11	278	15614	2619497	81	138
Spring 11/12	317	20200	5879742	177	359
Fall 11/12	338	21946	7167023	91	170

**Fig. 2** Number of Sessions per User (Hourly) (a) Moving Average, and (b) CDF

the evolution of the usage across the hotspot throughout the academic semesters.

In general, an increasing trend is observed in the number of deployed APs, number of wireless users and overall number of RADIUS sessions (start, alive, and stop), from semester to semester. Total input and output traffic, however, fluctuate between spring and fall semesters to some extent. Although the overall sent and received traffic grows in volume in ultimate fall/spring semester rather than the earlier, the wireless network are subjected to higher traffic in spring semesters compared to fall semesters.

4.2 Preliminary Data Analysis

In this section we present some extensive statistical analysis about the entire dataset and demonstrate relevant graphics revealing some general facts of underlying usage pattern of FEUP wireless network. We conduct this study from two peculiar viewpoint, users and the accompanying sessions, and APs and their accommodated users.

4.2.1 User Sessions

As indicated earlier, each user could connect to the same AP more than once during the day, and each connection creates a separate sessionID in the accounting table. An ideal association to the wireless network could last for the entire day and if the user is fixed in its lo-

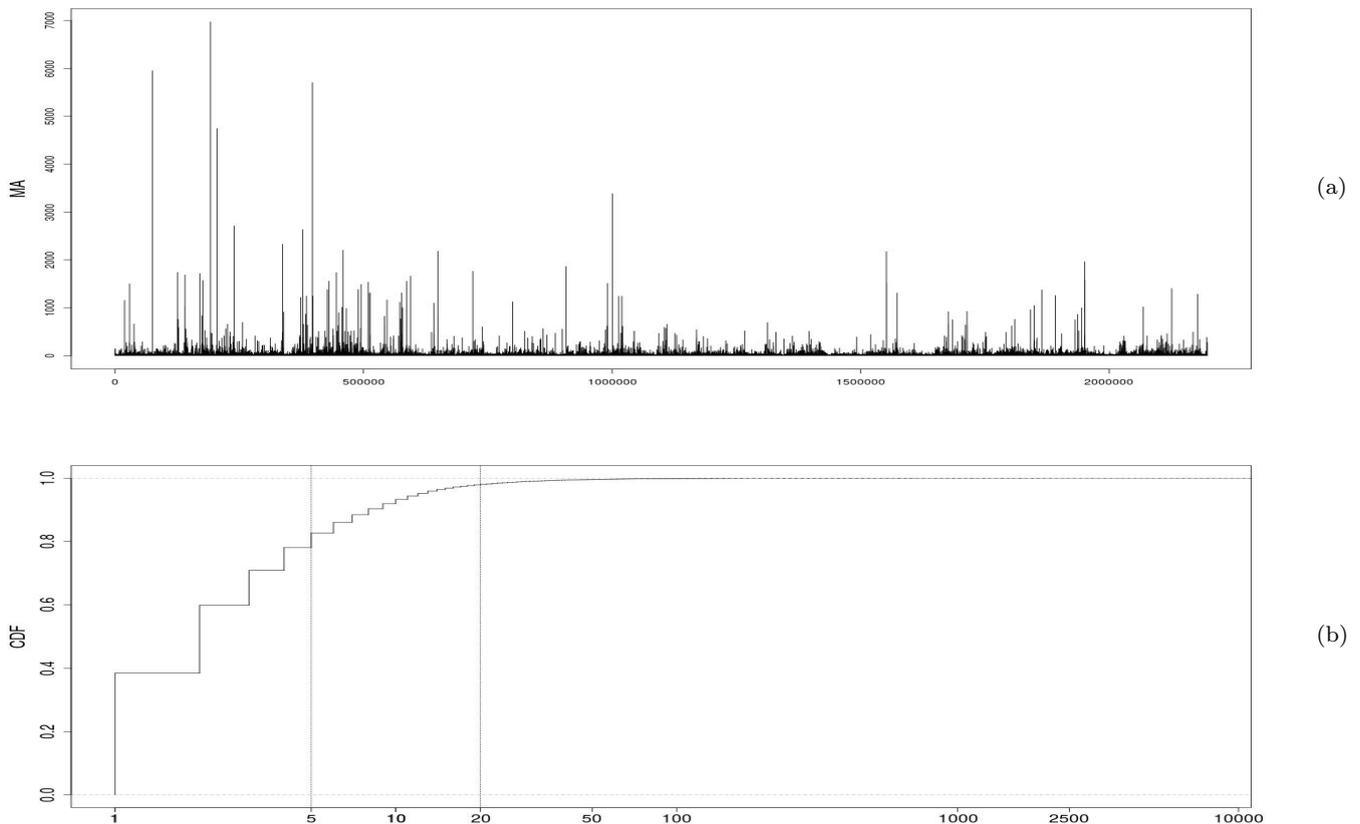


Fig. 3 Number of Sessions per User (Daily) (a) Moving Average, and (b) CDF

cation, it is expected to have the same session without interruption. However, this is not always the case and users disassociate from their current AP and associate to the same AP or another AP in the vicinity for various reasons. Figure 2 considers the proportion of the user sessions in an hourly period and Figure 3 reveals the same information on a daily basis.

Figure 2a shows the moving average of the number of sessions that each client (device) creates during one hour of connection. Although the majority of users have a few number of sessions in an hour which shows few number of disassociations, the extreme cases are also detectable in this figure. For instance, users are observed that generate over 2000 sessions on average in an hourly connection to a single AP. To study the greatest population of users, Cumulative Distribution Function (CDF) of users and their containing sessions is demonstrated in Figure 2b. This figure displays that more than 70% of the user connections remain unbroken and preserve a single session during the hourly association to their affiliated AP, and over 95% of the user connections contain only 5 sessions during an hour which is the result of intentional or unintentional disassociation from the current AP.

Figure 3a encloses similar information as Figure 2a, but in a daily basis. As expected, the number of disassociations during one day is higher than an hour period. Figure 3b demonstrates that extremely consistent connections which hold a single session during a day, are less than 40%. Most of such connections could be issued from stationary idle devices in vacant locations of the campus with few or no other active users around. This figure also displays that about 20% of the sessions are interrupted between 5 and 20 times a day.

4.2.2 Access Points

In this part, the study is more focused on the usage behavior of APs as indicators of different locations around the university campus. Figure 4 demonstrates the average number of users and sessions per AP during the two years of experiment for the working days only. Clearly this statistics could differ from semester to semester as the number of users and their corresponding sessions evolve over time, however this figure provides a general report of involvement of the entire set of APs in

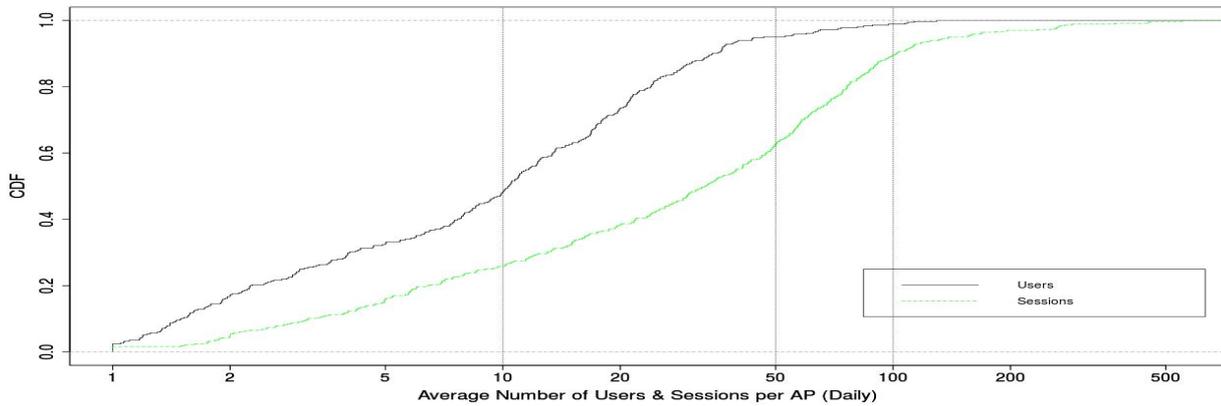


Fig. 4 CDF of Average Number of Users & Sessions per AP

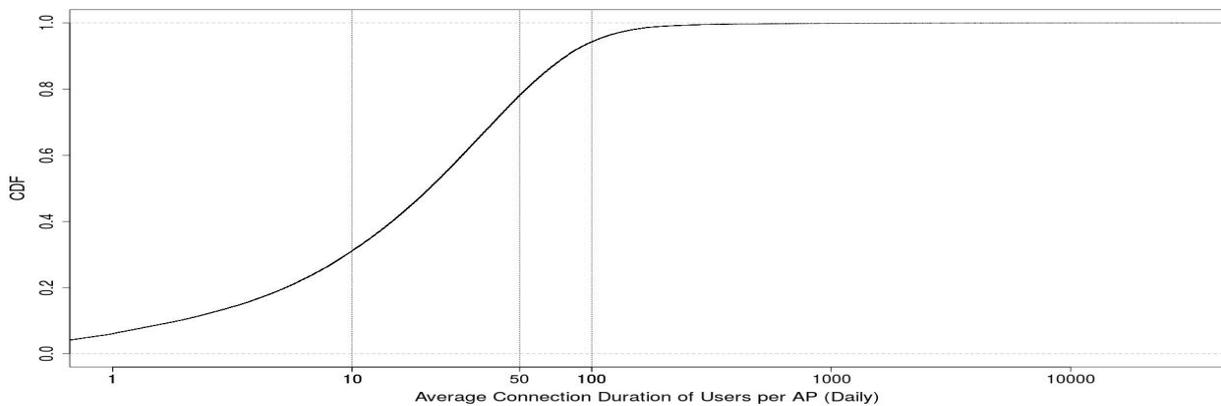


Fig. 5 CDF of the Daily Average Connection Duration of Users per AP (min)

the wireless covered area in two years of experiment. Figure 4 displays that around 20% of the APs contain only 10 sessions per day and almost 45% of the APs associate with 10 users during the day. To maintain the steadiness of the results, the weekends are excluded from this statistics. The figure also shows that over 95% of the APs (345 APs) contain at most 50 users a day and about 30% of the APs (109 APs) typically associate with 5 users each day.

Figure 5 reveals interesting information on the duration of users' daily connections per AP. It shows that the average connection period of users in 30% of the time is only 10 minutes per day. This data most probably belongs to the mobile users, guests, short-term clients or inactive users. Figure 5 also demonstrates that around 95% of the time, users maintain their connections to APs at most for 100 minutes, less than 2 hours a day. Such information send an important message to the network managers of the vitality of connection performance and quality of service as a great number of the users are connected to the network for less than 2

hours a day and getting interrupted over and over again in such a short period of time could be disappointing.

The study of the information provided in Figure 4 and 5- yields more precise understanding of the importance of the APs and learning their usage pattern based on the locations, for instance whether they are located in a busy entrance hall or a quiet corner of the campus. Such sort of information also imply the potential categories in terms of university divisions like administrative office, classroom, cafeteria, auditorium, etc. Such classification plays an important role for further analysis and modeling practices for the purpose of anomaly detection. It brings about the question of similarities (or differences) of the usage patterns in potential groups with different population of users that could prompt interesting anomaly detection strategies by learning the trend of the group and detecting the unusual events. These lines of research are of our interest for the future work.

4.3 Data Features

A number of features emerge from the raw dataset as a result of a preliminary analysis and enumeration process on a timely basis of 15 minutes. We categorize all the measured features as two main classes: *Density Attributes* and *Usage Attributes*. Those features that are indicators of *density*, basically demonstrate how crowded is the place in terms of active attendant users, when in fact the *usage* features disclose the volume of sent and received traffics by the present users. The former attributes mainly characterize the association population and durability, and the later ones reveal the total bandwidth throughput regardless of how populous is the place and it is more relevant to the applications utilized by the current mobile users.

4.3.1 Density Attributes

User Count : the number of unique users observed in a specific location (indicated by an AP) during the pre-defined time-slot (15 min).

Session Count : the total population of active sessions during a time-slot regardless of the owner user. This attribute reveals the number of attempts made by the the congregation of the present users to associate to the current AP. The connection time span of each user consists of one to many sessions.

Connection Duration : the total duration of association time of all the current users. This attribute is an indicator of the overall connection persistence. The utmost amount of this features is achieved when there is no evidence of disassociation in the ongoing active sessions during a time-slot ($User\ Count * 15\ min$).

4.3.2 Usage Attributes

Input Data in Octets : the number of octets transmitted from the client and incoming to the NAS port, and is only present in the Stop or Alive sessions. This attribute briefly refers to the number of bytes uploaded by the wireless user.

Output Data in Octets : the number of octets received by the client and leaving the NAS port, and is only present in the Stop or Alive sessions. This attribute shortly refers to the number of bytes downloaded by the wireless user.

Input Data in Packets : the number of packets transmitted by the client and incoming to the NAS port. This attribute is similar to the above *Input-Octet*, just to be measured in packets instead of bytes.

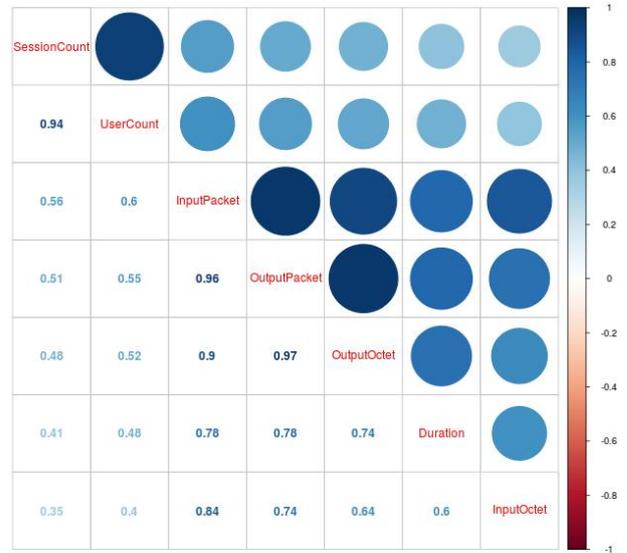


Fig. 6 Correlation Matrix of the Main Data Features

Output Data in packets : the number of packets received from the client and leaving the NAS port. This attribute is similar to the above *Output-Octet*, just to be measured in packets instead of bytes.

4.4 Feature Selection

In this section we discuss the connection and correlation of the data features explained earlier and disclose how to choose the best set of features for further analysis.

Figure 6 depicts the correlation matrix of all the above features. There is a high correlation observed between *User Count* and *Session Count*, on the grounds that the number of sessions are always equal or higher than the number of users in a time-slot. *Duration* do not have a strong correlation with any of the mentioned features, neither with *Density Attributes*, nor with *Usage Attributes*.

Having considered the input and output traffic transferred in octets, there is no significant correlation between these two compared to Input and output data in packets. However there is a slightly noticeable correlation between *Output Octets* and its corresponding attribute *Output Packets*, as well as *Input Octets* and *Input Packets*. Although there is a slight correlation between input/output data in octets and in packets, we consider them as semi-independent variables and include both of them in our further experiments. The in-

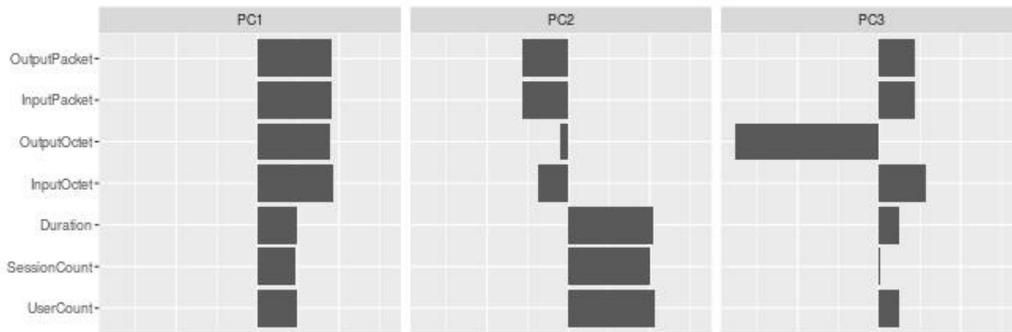


Fig. 7 The Behavior of the Main Features Relative to the Three Principal Components

formation added to the system through input and output traffic in octets simply take into account all the sent and received data in bytes. However, the input and output traffic measured in packets, could bring other types of information as the packets' size could differ by various factors such as application types and communication protocols.

For subsequent analysis and modeling procedures, we favor using less features rather than the entire set of attributes introduced earlier. For this reason, we applied Principal Component Analysis (PCA) technique to find the combination of the variables which best explain the phenomena and contain the greatest part of the entire information.

In this case the first three principal components bring the cumulative proportion of variance to over 95%. Figure 7 demonstrates the participation proportion of each feature to the principal components. From an analysis of this Figure we conclude that the first principal component is associated with all the above features in a positive manner, more specifically with the usage attributes. The second principal component is declined towards the density attributes, increasing with the larger density values and yet decreasing with the larger usage values. The single largest contributor to the third principal component is the input data in octet or the amount of uploaded bytes by the wireless users. The other features play less important roles in the third component, positively or negatively. Approximately categorizing the principal components like so, provides us with a deeper understanding of the connection of the aforementioned features, density or usage attributes, with the emanate best features resulted by PCA technique.

4.5 Conclusions

In this section we introduced collected RADIUS data from FEUP hotspot as the main dataset of this work

and performed a preliminary analysis from two points of view- user sessions and access point- to demonstrate the situation of the data respecting the hourly and daily sessions per user as well as the user population and connection duration per AP. Moreover we presented two main groups of features- usage attributes and density attributes- and defined a number of features for each group. We further studied the connection and correlation of the data features and selected the best combination of those features applying PCA. In the upcoming section we show how to model the AP usage data using the selected features represented in this section.

5 Statistical Modeling of 802.11 AP Usage

In this section we introduce statistical techniques for modeling purposes and in the upcoming section we indicate how to apply these models for anomaly detection. The modeling approach itself can be used in distinct directions such as to study the similarities and differences of the locations, to categorize the localities in terms of functionality (e.g. classroom, office, library) or specification (homogeneous/heterogeneous daily, seasonal or constant usage). We introduce time-invariant and time-variant models and in each case we show how to apply the model on the large dataset previously elaborated.

5.1 Time-invariant Modeling

We first consider models that assume there is no time binding between consecutive daily events. Although this might not be precisely the case, it yields simpler modeling approach. Later in the paper we compare this type of modeling with others that do consider dependency between consecutive daily events.

5.1.1 Gaussian Mixture Model

We begin our modeling efforts by applying techniques that assume all daily events come from the same distribution, regardless of any time dependency between the consecutive records. To explain this, we pick Gaussian Mixture Model (GMM), a probabilistic model that presume all the data points are generated from a mixture of a finite number of Gaussian distribution with unknown parameters. The Expectation Maximization (EM) procedure is the optimization technique utilized to fit the unknown parameters and incorporate information about the covariance structure of the data as well as the centers of the latent Gaussians [28].

$$p(x|\lambda) = \sum_{k=1}^M \omega_k g(x|\mu_k, \Sigma_k) \quad (1)$$

where x is a D-dimensional continuous-valued data vector (of features), $w_k, k = 1, \dots, M$, are the mixture weights, and $g(x|\mu_k, \Sigma_k), k = 1, \dots, M$, are the component Gaussian densities. Each component density is a D-variate Gaussian function of the following form,

$$g(x|\mu_k, \Sigma_k) = \frac{\exp\{-\frac{1}{2}(x - \mu_k)' \Sigma_k^{-1} (x - \mu_k)\}}{(2\pi)^{D/2} |\Sigma_k|^{1/2}} \quad (2)$$

with mean vector μ_k and covariance matrix Σ_k . The mixture weights satisfy the constraint that $\sum_{k=1}^M \omega_k = 1$.

The complete Gaussian mixture model is parameterized by the mean vectors, covariance matrices and mixture weights from all component densities. These parameters are collectively represented by the following notation,

$$\lambda = \{\omega_k, \mu_k, \Sigma_k\} \quad k = 1, \dots, M \quad (3)$$

5.1.2 GMM Application: Case Study

GMM could be applied to our data features in several ways, for instance a single mixture model for the entire set of data, or a mixture model for each location separately. The later approach is closer to our goal of proposing practical models for each place indicated by an AP (or a broader neighborhood) to explore the characteristics of that place, and ultimately discovering the abnormal behaviors occurring in contrast with the expected usage pattern. Note that in our previous work [10] we modeled and identified the anomalies of three categories: of a single model for all APs, a mixture model for groups of APs and individual models for each AP. In this work we study the individual model to be able to evaluate it with our deployed testbed and

in the future work we intend to explore the models for the potential groups of APs.

In order to investigate the modeling capacities of GMM for the mentioned aims, we select two different spots to be our test cases: a highly crowded AP at the computer service section with 3726 observed users, and a less crowded AP in the chemical engineering department with overall 175 users. The experiment takes into consideration the second semester period of 2011 from February to July. To achieve more precise result, we focus on the working daily pattern, hence the data records belong to the working days (from Monday to Friday) and the working hours (8 a.m. to 6 p.m.).

On each location, GMM fits are computed with three mixture components. The Gaussian density parameters (mean and covariance matrix) are depicted in Figure 8, the first row belongs to the crowded AP and the second row shows the density parameters of the less crowded AP. In order to facilitate the visual perception and to have an easier comparison, the density parameters are illustrated in 2D, despite the fact that GMM process is conducted on 3 features (principal components).

The data is standardized on each column to have zero mean and one standard deviation, so the density values are not appropriate to be compared with each other directly. However, the contour lines show the diversity of the data points in each mixture component and the direction of spread as well as the mass center. The R value on each plot represents the correlation between the X and Y axis, correspondingly the first two principal components.

Each location is characterized in this manner and according to GMM modeling approach,

$$\lambda_1 = \{\omega_{i1}, \mu_{i1}, \Sigma_{i1}\} \quad i = 1, \dots, 3$$

and

$$\lambda_2 = \{\omega_{j2}, \mu_{j2}, \Sigma_{j2}\} \quad j = 1, \dots, 3$$

represent the mixture weights and density parameters of the first and the second APs respectively.

5.2 Time-variant Modeling

In this section we consider models that assume time dependency between consecutive daily events. In this case the sequences of data records matter and they form significant connections in a meaningful context or profile. In time-variant models in general, conditional probabilities for events are determined based on the history of the events. In the following section we study the Hidden Markov Models for modeling the time-varying sequential data for the ultimate purpose of anomalous pattern recognition which we discuss more in detail in the next section.

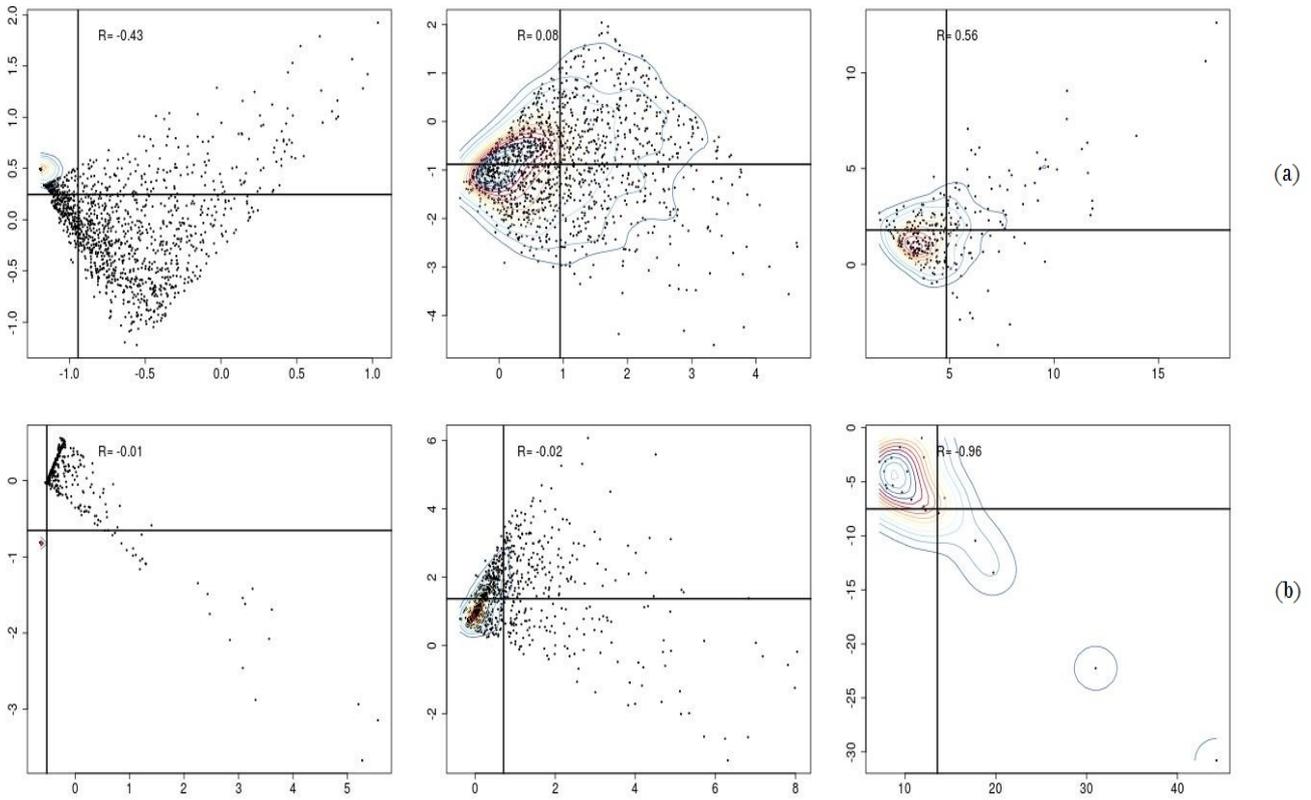


Fig. 8 Density Parameters of Three Gaussian Mixture Components of the Selected APs. **a** Crowded AP, **b** Less Crowded AP

5.2.1 Hidden Markov Model

HMMs are generally used for the stochastic modeling of non-stationary time-series. HMMs provide a high level of flexibility for modeling and analyzing time-varying processes or sequential data. Their particular application is in recognition such as speech recognition, activity recognition, gene prediction, etc. where data instances are represented as a timely sequence of estimates. In the current research we propose how to use HMMs for modeling and anomaly detection purposes in wireless networks which has never been investigated before to the best of our knowledge.

Rabiner and Juang [27] presented a comprehensive tutorial on HMM which provides a profound understanding of the basic blocks of HMM. HMM symbolizes a doubly stochastic process with a set of observable states and a series of hidden states which can only be observed through the observable set of stochastic process. The goal in HMM is recovering a data sequence that is not immediately observable through the other set of observable data.

The formal definition of a n -state HMM notation is determined as follows:

- A set of hidden states $S = \{s_i\}$, $1 \leq i \leq n$
- A set of possible symbol observations in discrete models $V = \{v_i\}$, $1 \leq i \leq m$
- State transition probability distribution (transition matrix) $A = \{a_{i,j}\}$, $1 \leq i, j \leq n$, $a_{i,j} = P(s_j \text{ at } t+1 | s_i \text{ at } t)$
- Observation symbol probability distributions (emission matrix), $B = \{b_j(k)\}$, $1 \leq k \leq m$, $b_j(k) = P(v_k \text{ at } t | s_j \text{ at } t)$
- Initial state distribution $\pi = \{\pi_i\}$, $1 \leq i \leq n$, $\pi_i = P(s_i \text{ at } t = 1)$
- $m =$ number of observation symbols in discrete models
- $n =$ number of hidden states

The set $\lambda = (A, B, \pi)$ completely defines an HMM [27]. However, in continuous emissions, instead of having m outcomes for the observations, distribution parameters such as mean and covariance are determined. In such cases a model is represented as $\lambda = (A, \mu, \Sigma, \pi)$, and μ and Σ stand for mean vector and covariance matrix respectively.

Using the model λ , an observation sequence $O = o_1, o_2, \dots, o_T$ is generated as follows:

1. Select an initial state, s_1 , according to the initial state probability distribution, π ;

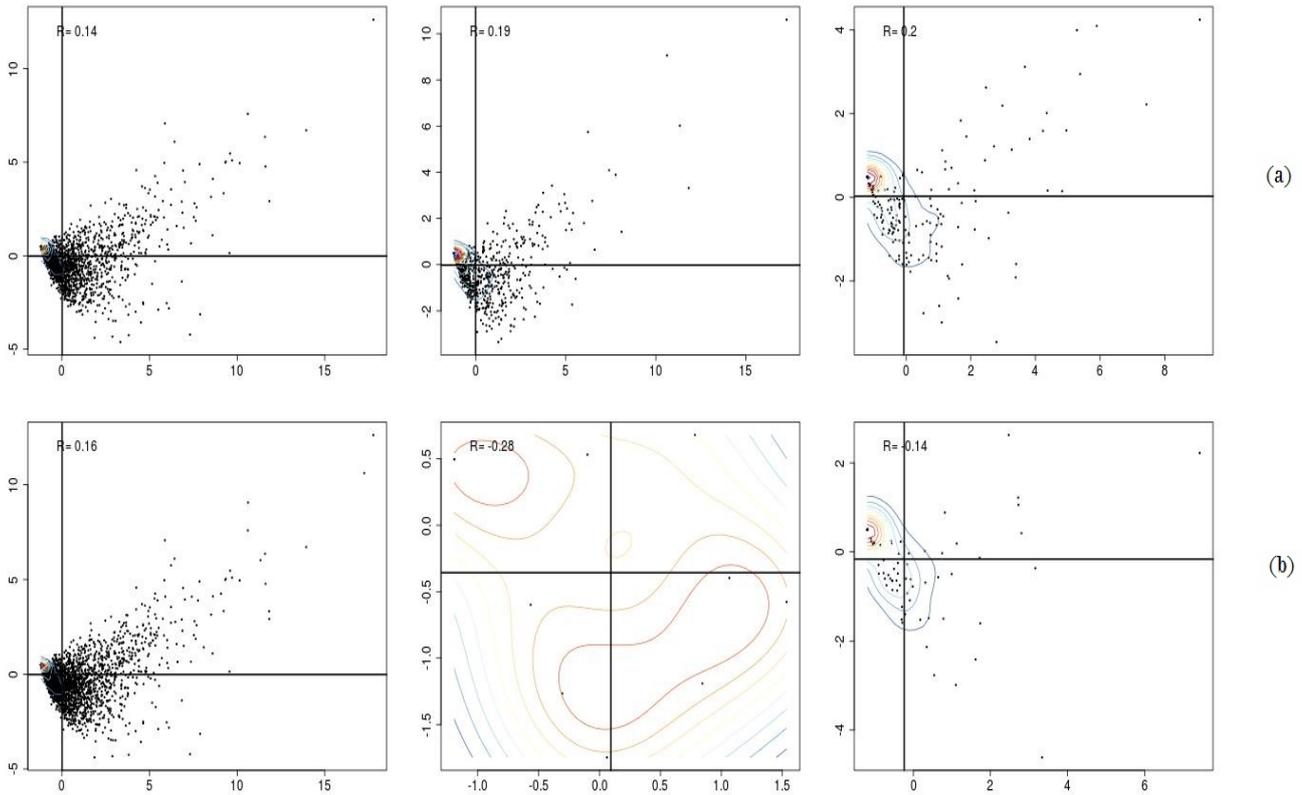


Fig. 9 Density Parameters of Three Hidden State in HMM of the Selected APs. **a** Crowded AP, **b** Less Crowded AP

2. Set $t = 1$;
3. Choose o_t according to observation probability distribution in state s_t , $b_{s_t}(k)$;
4. Choose s_{t+1} according to the state transition probability distribution for state s_t , $a_{s_t, s_{t+1}}$
5. Set $t = t + 1$; return to step 3 and continue until $t > T$

Given the form of the HMM, there are three key problems of interest that solving them promotes modeling the real world applications. These problems are listed as the following [27]:

Problem 1 – Given the observation sequence $O = o_1, o_2, \dots, o_T$ and the model $\lambda = (A, B, \pi)$, how we compute $P(O|\lambda)$, the probability of the observation sequence.

Problem 2 – Given the observation sequence $O = o_1, o_2, \dots, o_T$, how we choose a state sequence $S = s_1, s_2, \dots, s_T$ which is optimal in some meaningful sense.

Problem 3 – How we adjust the model parameters $\lambda = (A, B, \pi)$ to maximize $P(O|\lambda)$.

According to our data set, the HMMs form observations with continuous multivariate Gaussian distribution, hence the emission matrix B is defined by the distribution parameters associated with the set of states. In the proposed model, the HMMs contain fully con-

nected states, thus transitions are allowed from any state to any other state.

5.2.2 HMM Application: Case Study

In this section we select the very same APs as in the GMM case study (Section 5.1.2), and build HMM models for each of them separately. Our focus is once more on the working daily pattern in the second semester of 2011, from Monday to Friday in the working hours.

As described earlier, we consider fully connected HMMs (ergodic model) with continuous Gaussian distribution as the emission probabilities and 3 hidden states. The states are initialized randomly, and the number of states is selected heuristically based on the best practice of the experiments conducted on both the large dataset and the Testbed dataset. For the multivariate Gaussian distribution of the observations, each component of the mean vector is uniformly drawn between $\mu - 3\sigma$ and $\mu + 3\sigma$ and the initial covariance matrix is diagonal and each initial variance is uniformly drawn between $\frac{1}{2}\sigma^2$ and $3\sigma^2$. The initial probability matrix (π) and the transition matrix (A) are uniformly drawn. The initial HMM is then optimized with the Baum-Welch algorithm with the cut off likelihood value of

$1e - 6$ or the maximum number of iterations set to 20. After the optimization process, the physical meaning of the hidden states are more discernible. The values of the principal components in each state shows the tendency of the states to the usage or density attributes. For example a hidden state with the highest value for the second principal component shows a more populated case in terms of users or sessions density. In the future work where the concern is more on modeling the anomalous patterns we utilize the interpretation of the hidden states to relate them to the physical conditions of the locations.

The Gaussian density parameters of the three hidden states are illustrated in Figure 9, similar to Figure 8, the first row is affiliated with the crowded AP and the second row belongs to the less crowded AP. The contour lines in these two figures represent the overall picture of the population and density distribution of the data in each component or state. Suchlike graphs are visual aids to depict the density parameters only, and for inspecting the goodness of distribution over the entire feature set and make any comparison, further investigations are required.

5.3 Model Comparison: GMM vs. HMM

In this section two techniques are considered only for the sake of modeling purposes, a time-invariant model (GMM) and a time-variant model (HMM). In the coming section we investigate the ultimate goal of this modeling which is the recognition of anomalous points or regions. At this stage, before exploring the anomaly detection territory, we briefly itemize the modeling functionalities and propose some simple tests to verify the more qualified model.

The potential functionalities of the locations characterization and modeling are listed as following:

- Classification of the locations, represented by APs, in terms of utility and temporal patterns.
- Recognition of the meaningful similarities and distinction of the locations.
- Grouping the most related APs and propose mixture models for the groups [10].

To investigate the competency of the two proposed models and estimate the capacity of each, we conduct a simple test. First of all, we measure the log-likelihood of the models in modeling the training data of the two samples, crowded AP and less crowded AP, and then we select a random day from each AP and calculate the log-likelihood of the models towards the test data which is new to both models. We use log-likelihood values (LLV) to measure the goodness of fit of our models.

The model with larger log-likelihood value surpasses the model with smaller log-likelihood value.

Given data x with independent multivariate observations x_1, \dots, x_n , the likelihood of a Gaussian mixture model with M components is defined as [14]:

$$likelihood(x|\lambda) = \prod_{i=1}^n \sum_{k=1}^M \omega_k g(x_i|\mu_k, \Sigma_k) \quad (4)$$

where $g(x|\mu_k, \Sigma_k)$ is the k th component's Gaussian density, as already defined in Equation 1, and ω_k is the probability that an observation belongs to the k th component.

The log-likelihood function takes the following form:

$$\log\text{-likelihood}(x|\lambda) = \sum_{i=1}^n \log\left(\sum_{k=1}^M \omega_k g(x_i|\mu_k, \Sigma_k)\right) \quad (5)$$

In the EM process, the parameters of the GMM, λ , are estimated so that the likelihood of the GMM given the training data is maximized, Maximum Likelihood Estimation (MLE). Ensuing several iterations, the MLE yields the likelihood of the GMM given the training data. We applied MClust R package [15] to conform the Gaussian mixture components and estimate the log-likelihood of the training and test data provided in Table 3.

The likelihood of a HMM is basically the first key problem of HMMs stated earlier, the probability of an observation sequence given the model parameters:

$$\begin{aligned} P(O|\lambda) &= \sum_{all S} P(O|S, \lambda)P(S|\lambda) \\ &= \sum_{s_1, s_2, \dots, s_T} \pi_{s_1} b_{s_1}(O_1) a_{s_1, s_2} b_{s_2}(O_2) \dots a_{s_{T-1}, s_T} b_{s_T}(O_T) \end{aligned} \quad (6)$$

We utilized GHMM library [29] for the formation of HMMs, estimation of log-likelihoods and all the other requisites of the experiments performed in this work.

Table 3 contains the log-likelihood values of the trained GMM and HMM models for the selected APs, regarding both the training and test data. Comparing the log-likelihood values of the training data, HMM provides higher values (less negative) both for the crowded AP and the less crowded AP. Note that the training data contains 25 days data and the test data consists of only one day data selected randomly from the unobserved days. Concerning the test data, it is expected that the selected day from the same AP obtains higher log-likelihood value rather than the data from another AP due to the possible similarity of daily usage in a specified location. The first GMM (built over the crowded AP data) provides the same amount of log-likelihood for both test data, thus yields no distinction for its own usage pattern rather than the other AP. However, the second GMM (trained with the less crowded AP

Table 3 Log-likelihood Values (LLVs) of the Training and Test Data Belong to the Selected APs for GMM and HMM Models

Test Data LLVs	Trained Model	GMM	GMM	HMM	HMM
		Crowded AP	Less Crowded AP	Crowded AP	Less Crowded AP
The same train data		-3468	-2154	-2553	-2131
Test data from the crowded AP		-189	-189	-134	-209
Test data from the less crowded AP		-509	-95	-195	-115

data) provides higher log-likelihood value for its own data rather than the other AP.

HMMs, on the other hand, provide higher amount of log-likelihood for their own test data rather than the other AP, which shows the better matched model for self data. Both GMM and HMM models for the crowded AP provide close values of log-likelihood for the test data, so the models do not seem to be very robust in distinguishing between its own data and the other AP. However, GMM and HMM models for the less crowded AP achieve higher log-likelihood values for the training data rather than the models of the crowded AP. It must be considered that the test data is selected randomly and the pattern of the selected day is not determined in terms of normal or abnormal usage, nevertheless the overall outcome of HMM models looks more satisfying compared with GMM. In Section 7, the experiments are conducted on the testbed dataset with recognized anomalies so that the conclusion will be based on the known ground truth. In the next section, we investigate the time-variant specifications of HMMs towards the simplicity of the time-independent GMM concerning the anomaly detection objectives.

5.4 Conclusions

In this section we presented GMM as time-invariant and HMM as time-variant modeling techniques. As a case study for each approach we selected two different locations in the university campus- a highly crowded AP and a less crowded AP- and applied the forenamed methodologies. We then defined the log-likelihood for each method separately to examine the goodness of fit for the proposed models in terms of train and test data. Having conducted a simple experiment on the selected APs revealed that HMMs are more likely to provide a robust model to distinguish between their own pattern and an unfamiliar pattern. In the following section we show the functionality of the proposed models to detect anomalous cases in AP usage data.

6 Detection of Anomalies in AP Usage Data

In this section we show how the aforementioned models are utilized for the purpose of anomaly detection. We further explore the capabilities of these models in recognition of abnormal events and series of unexpected occurrences.

6.1 Anomaly Detection Approach

Network administrators are generally concerned with anomaly detection as well as prediction. These two important tasks enable them not only to make immediate decisions to alleviate the complications of the network, but also to establish longstanding plans to support the expansion of the network and its dynamic usage over time.

6.1.1 GMM Estimation: Divergence from the Gaussian Densities

The most generic definition of the anomalies asserts those points or small regions isolated from the normal zones which contain the majority of the observations. Thus, a straightforward approach to detect anomalies, when there is no ground truth available, is to define the normal zones and distinguish those rare observations which hardly belong to those normal sectors.

In GMM, the time-invariant model discussed earlier, a number of Gaussian mixture components are determined and each component contains normal density parameters. The model is built based on several training data and the newly arrived records are inclined to the most compatible component with the least distance. Hence, to detect abnormal points we need to estimate the affinity degree of each point, as already described in Equation 4, and mark outliers as having the slightest probability of belonging to any cluster.

6.1.2 HMM Estimation: Likelihood Series

HMM, as a time-variant model, considers the temporal dependency between consecutive data records. Calculating the log-likelihood of a single data point or a series

of sequential data points as already expressed in Equation 6, emanates the mis-behaving records comparing to the log-likelihoods of the norm of the data. The unexpected low values for the log-likelihood in HMM are generally due to one or some of the following arguments:

Divergence from the Assigned Hidden State: Given an HMM model λ and an observation sequence of $O = o_1, o_2, \dots, o_T$, the most probable set of states are generated by *Viterbi* algorithm as $S = s_1, s_2, \dots, s_T, s_i \in S$. To estimate the distance of a data point in time t to its counterpart HMM state (s_t) in *Viterbi* path, the *Mahalanobis* distance is evaluated between time-series elements and the hidden states. Consequently the outliers which display the unreasonable distance to their assigned hidden states, are highlighted to potentially have a poor value in the likelihood series. This approach is approximately similar to the outlier detection technique addressed earlier for GMM components.

Less Likely State Transition: According to the third well-known HMM problem, stated in Section 5.2.1, when a HMM model is trained the model parameters are adjusted to maximize the probability of the observed data $P(O|\lambda)$. The transition probability matrix is one of the salient components of the trained model. The highest transition probabilities are frequently observed between identical states (s_i to s_i), while the lowest probabilities often occur between the most distant states. However, regardless of the form of the transition matrix, in the *Viterbi* sequence outcome, it is expected to observe the transition probabilities proportional to the values of the transition matrix. Whenever this principal is violated there exist an anomaly prospect. For instance if in a *Viterbi* path the transition from state s_i to state s_j occurs only once (out of 60 transitions), and the transition probability of $a_{i,j}$ is 30% in the transition matrix, this circumstance sounds unlikely and thus an anomaly-prone transition.

6.1.3 Anomaly Detection: Case Study

In this section we explore the addressed methodologies to detect anomalous data points or data sequences in the same two APs that we proposed GMM and HMM models for their corresponding training data. Figure 10 highlights the outliers of the one day test data detected by measuring the largest distance from the Gaussian components. The result of the first AP (crowded AP) is displayed in blue and the second AP (less crowded AP) is demonstrated in green. Two data points are marked in red that both belong to the first model of the crowded

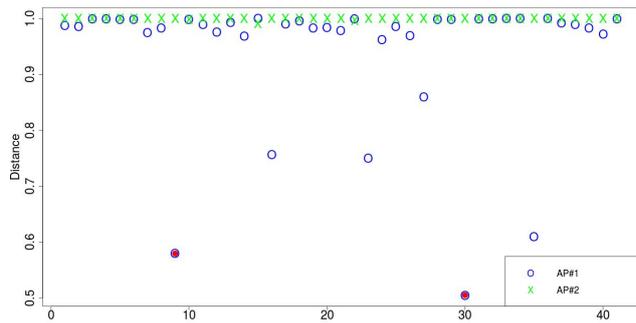


Fig. 10 GMM Estimation of Anomalous Data Points Based on the Largest Distance from the Assigned Gaussian Component

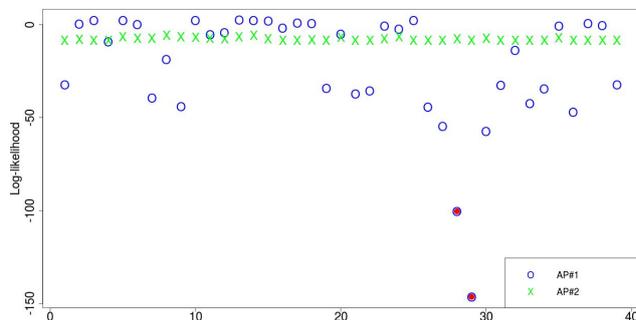


Fig. 11 HMM Estimation of Anomalous Data Points Based on the Lowest Log-likelihood

AP. These outliers are appointed to a Gaussian component of the first model, but with the lowest probability (less than 60%). Here we selected the normality threshold to be 60%, however it could differ from model to model and the most appropriate value of threshold could eventually be decided by the network manager.

Figure 11 displays the anomalous points detected by HMM based on the lowest value of the log-likelihood. In this approach, two different data points are marked as outliers which belong to the first AP training data, the crowded AP. The cut-off value is considered to be log-likelihoods below -100, note that this value could also be configured. The more strict cut-off value yields higher false positive rate. We investigated the likely origins of the outliers emerged in this case and we observed that the *Mahalanobis* distance of the marked data points are maximal with the assigned hidden state in the *Viterbi* path. That must have caused the low log-likelihood value in the likelihood series. Further experiments on anomaly detection by HMMs and evaluation techniques are performed in our previous work in [10].

However, in this case study we demonstrated how the anomaly detection analysis work in our proposed framework. In the next section, we evaluate both models based on the achieved results of the deployed testbed, hence we can determine with more confidence which points are detected correctly.

6.2 Conclusions

In this section we described the anomaly detection techniques by GMM and HMM. In GMM we define anomalies as the distant data points that hardly belong to any Gaussian component, while in HMM anomalies are the data points with the minimum likelihood value. As discussed more in detail in our previous work [10], we analyzed the root cause of the low likelihood value as divergence from the assigned hidden states as well as the low probability in state transition. We further explored the addressed methodologies to detect anomalies at the same APs from the previous section. We justified the detected anomalous points, however in absence of the ground truth in the large dataset it was not possible to thoroughly evaluate the anomalous points and we left the evaluation process for the next section.

7 Experimental Setup

In order to validate anomaly detection techniques proposed in this work we deployed an exploratory testbed with one single AP and generate a number of anomalies in a controlled environment for experimental purposes. We work with FreeRADIUS server which is widely used for Enterprise Wi-Fi and IEEE 802.1X network security and communication, particularly in the academic community, including Eduroam [1]. The very basic aspects of our testbed dataset is elaborated in the following section.

7.1 Server Configurations and Users Specifications

The set up process of the FreeRADIUS server is performed on a Linux machine with 2.30 GHz Intel(R) Core(TM) i5-2410M CPU, and 8GiB System Memory. The database system used to store primary configurations and AAA information is MySQL and consist of 10 preordained tables. The principal tables employed for data collection and analysis are labeled as radcheck (authentication), radpostauth (authorization) and radacct (accounting). Other essential configurations are conducted directly on FreeRADIUS setting files, such as

server and client security configurations, required certificates, database setups, and so forth.

As stated earlier, the testbed deployed for this study is dedicated to one-AP-many-users. Thus, we describe the AP configurations and wireless users specifications in the following lines. The AP is an enhanced 802.11g wireless access point powered by D-Link 108G technology, DWL-2100AP, and supports WPA and WPA2 security protocols. The wireless users connected to this network during one month of experiment consist of two laptops, two smart phones, and two tablets. A summary of the users' specifications in terms of devices, operating systems and participation time in the experiment is provided in Table 4. Obviously not all the users were present everyday and every hour of the test, but they follow a natural form of entering and exiting the network. Some devices were disassociated from the network when the users simply depart from the coverage area and others were deliberately disconnected in the time of specific anomaly generation. In the coming sections we present all types of anomalies generated and organized for this testbed.

Table 4 A summary of the testbed users' specifications

Device	OS	Participation Time (%)
Surface Pro II	Win 10	100%
Asus	Win XP	100%
Alcatel onetouch	Andriod	85%
iPhone	iOS	15%
iPad	iOS	100%
Dell	Win 10	8%

7.2 Network Anomaly Generation in a Controlled Environment

In this section we describe how some of the known wireless network issues are re-generated to make the desired data records for the evaluation of the proposed methods in this work. In the time of experiment not all days encounters anomalies, some days simply end as NORMAL days and the users' connection and amount of network usage are according to the users' usual plan of the day. In ABNORMAL days, however, one or some kind of anomalies are provoked to test the behavior of the model under abnormal circumstances. The anomalous patterns selected for this purpose are common cases occur in real networks relatively often and affect the performance of users connection and availability of the network. Succeeding paragraphs deal with the specific

aspects of these anomalies and point out how to replicate them.

7.2.1 AP Shutdown/Halt

To reproduce this anomalous effect, when there is no session recorded in the accounting table, the AP could be shutdown for a while or restarted. This anomaly is regenerated under various circumstances and for different period of time and in the real world could be considered as AP shutdown, halt, crash or power off.

7.2.2 Heavy Usage

Single User This anomaly arises when only one user performs heavy download or upload. It might affect the rest of the associated users depending on the amount of usage, duration, time of the day and other relevant factors.

Multiple Users This anomaly emerges when more than one user use the network excessively, and therefore the overall throughput of the network intensifies. This could occur in a NORMAL day or as an anomalous event and the network tolerance, as expected, varies for different networks and different AP configurations. In any case, the proposed model is expected to detect the irregularity and report the level of hazard so that the network managers could take control of the situation and make required changes if possible.

7.2.3 Wireless Network Interference

In a real network, a variety of things can interfere with the radio waves, degrading the quality of connection and decreasing the network reliability. Sources of interference are commonly from other wireless networks in the vicinity when they all locate in the same channel, from non-802.11 devices such as microwave ovens or cordless phones that use 2.4GHz band as well, from other clients in a crowded environment when they all try to transfer data at the same time, and from RF effects such as hidden terminals or capture effects. In this work we intend to cause interference anomaly in a systematic and controlled manner. For this aim, we made use of a python script named wifijammer [6] to intentionally jam wireless clients or APs in the range to simulate the same outcome as the aforementioned interferences. The jamming process works by sending 1 de-authentication packet to the client from the AP, 1 de-auth to the AP from the client, and 1 de-auth to the AP destined for the broadcast address to de-authenticate all clients connected to the AP. Many APs,

however, ignore de-auth to broadcast addresses. We employed wifijammer in the following plans by applying peculiar properties each time to create different forms of interferences.

Jamming the Entire Channel In this practice, the monitor mode interface is set to listen and de-authenticate clients or APs on a specific channel. This way of jamming influence all the available networks on the current channel and imply interferences caused by busy channels.

Jamming Clients with Various Time Intervals Executing the De-authentication procedure with short time intervals hinder clients from recovering and disable them for the entire period of jamming, so the immediate result in the accounting table is the one-time stop session from each client and then a silent period without any start session. While de-authenticating with a larger time interval makes clients reclaim and try to get back the connection to the AP, and subsequently many short sessions is observed in the accounting table because they are de-authenticated right after getting connected again. In such manner we can replicate two interference cases observed in the real datasets frequently.

Jamming Specific Clients De-authenticating some specific clients and not the rest, resembles the hidden-terminal situation, when one client is forced to back-off and delay data transfer because the other clients can not sense its send-request. Depending on the time interval discussed earlier, the sessions outcome in the accounting table could be different.

7.3 Testbed Experimental Results

The testbed experiment is deployed in a home environment, with a single AP and 6 regular users and between 3-4 guest users. The experiment contains 5 weeks of data, 30 working days, and is performed in two different time span, once in November 2015 and a while later in April 2016. There exist 20 normal days with no anomalies provoked, and 10 abnormal days containing at least one anomalous event a day. Each anomaly takes from 15 minutes to around an hour.

In the following paragraphs we show how the modeling and anomaly detection techniques operate in the presence of the ground truth, data obtained from the testbed deployment.

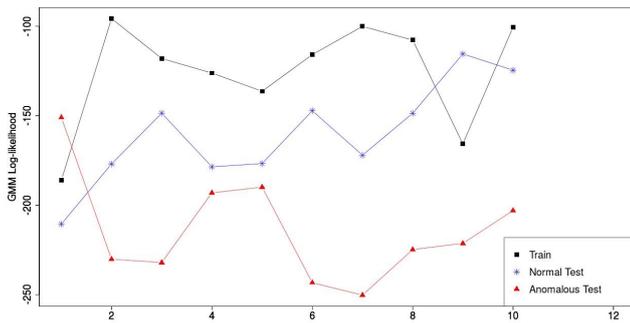


Fig. 12 Likelihood values of the training and test data belong to Testbed for GMM Model

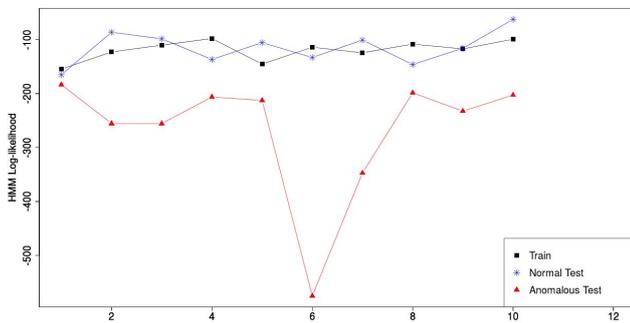


Fig. 13 Likelihood values of the training and test data belong to Testbed for HMM model

7.3.1 GMM vs. HMM Modeling: Pros and Cons

For the first experiment, a GMM model is built with 10 randomly selected normal days as training data. From then on, the likelihood of the generated model is computed against the training data as well as 10 unobserved normal days and 10 abnormal days as test data. The same process is performed on the HMM model, with the same set of training and test data. The summary of this experiment is displayed in Figure 12 and 13.

Both figures demonstrate overall higher likelihood values for the training data. The likelihood values of the unobserved data set is divided into normal and abnormal outputs which are displayed in graphs with different colors and shapes. In both models there are higher likelihood values for the normal days rather than the abnormal days. However, there is a discernible boundary between the normal and abnormal results in HMM while in GMM the likelihood values are not clearly separated and there are even some instances that the likelihood value of the normal day is lower than the abnormal day. The daily likelihood of abnormal days are apparently lower than the normal days, and this value

varies with the number of abnormal occurrences and duration of each event. However, it is more straightforward to define a threshold for HMM rather than GMM model, to announce a day normal or abnormal.

7.3.2 Anomaly Detection

In this section we determine the anomalous time-slots with the proposed methodologies and compare the achieved results from the model with the testbed anomalous ranges recorded for the abnormal instances. Note that various thresholds for each technique produce different results as the detection and false positive rates change based on the selected threshold. We made use of some statistical metrics to measure the detection accuracy and false alarms such as fall-out or false positive rate (FPR), specificity (SPC) or true negative rate (TNR), sensitivity or true positive rate (TPR), and eventually accuracy (ACC) and F1 score. In order to acquire the specific definition of each terminology refer to [5].

The summary of the analysis on the normal and anomalous test data are presented in Table 5 for GMM modeling and in Table 6 for HMM modeling approaches.

Table 5 shows that higher thresholds increase the possibility of anomaly detection (24.9% rather than 4.7%), however the false positive rates also increase accordingly (19% rather than 9.9% and 3%). In normal test data, when we expect no anomalies to occur, from 2.5% to 10.5% fall-out is observed. Comparing this fall-out ratio to the results of Table 6 for normal test set, it is noted that much lower false alarms is marked for HMM (from 0.5% to 3.75%). Furthermore, the FPR for the anomalous data in HMM is quite trivial relative to GMM FPR output (1.1% in HMM vs. 19% in GMM). The highest detection rate or TPR in HMM modeling is achieved with Threshold equals to -10 which is 75% in average for 10 abnormal days of the experiment.

Regarding the FPR or fall-out ratio recorded for normal data in HMM, a careful consideration on each false alarm is performed and it is noted that the HMM model is slightly sensitive to extreme download ratio and in some cases both download and upload volumes. As the testbed is deployed in a real home environment with real wireless users, although in normal days no anomaly is generated deliberately, there might have been some evidences of rather high download or upload by the users as it happens quite often in every wireless network. Therefore the false positive examples occurred in normal days could be introduced as real anomalies appearing in normal days, however for this experiment we assumed that normal days contain no anomalies. In our future work we intend to propose an unsupervised anomaly detection algorithm that detect anomalies in

Table 5 Anomaly detection of the normal and anomalous test data belong to Testbed for GMM

Statistical Metrics Data - Threshold	False Positive Rate (FPR)	True Negative Rate (TNR)	True Positive Rate (TPR)	Accuracy (ACC)	F1 Score
Normal Testset (Threshold: 0.6)	2.5%	97.5%	0%	97.5%	0%
Normal Testset (Threshold: 0.7)	5.5%	94.5%	0%	94.5%	0%
Normal Testset (Threshold: 0.8)	10.5%	89.5%	0%	89.5%	0%
Anomalous Testset (Threshold: 0.6)	3%	97%	4.7%	81%	8.1%
Anomalous Testset (Threshold: 0.7)	9.9%	1.01%	4.7%	75%	7.2%
Anomalous Testset (Threshold: 0.8)	19%	81%	24.9%	70%	20.75%

Table 6 Anomaly detection of the normal and anomalous test data belong to Testbed for HMM

Statistical Metrics Data - Threshold	False Positive Rate (FPR)	True Negative Rate (TNR)	True Positive Rate (TPR)	Accuracy (ACC)	F1 Score
Normal Testset (Threshold: -50)	0.5%	99.5%	0%	99.5%	0%
Normal Testset (Threshold: -20)	1.75%	98.25%	0%	98%	0%
Normal Testset (Threshold: -10)	3.75%	96.25%	0%	96%	0%
Anomalous Testset (Threshold: -50)	0%	100%	39%	90%	49%
Anomalous Testset (Threshold: -20)	0%	100%	43%	91%	52%
Anomalous Testset (Threshold: -10)	1.1%	98.9%	75%	95%	74%

Table 7 Detection rate of various anomalous patterns of the Testbed

Anomalous Patterns Model	Jamming Channel (Low Intervals)	Jamming Channel (High Intervals)	Heavy Usage (Single User)	Heavy Usage (Multiple Users)	AP Power Off
GMM (Threshold: 0.8)	28.5% (4/14)	17.3% (4/23)	8.3% (1/12)	0% (0/3)	35.2% (6/17)
HMM (Threshold: -10)	71.4% (10/14)	73.9% (17/23)	83.3% (10/12)	100% (3/3)	82.3% (14/17)

an unlabeled test dataset which is the case when no ground truth is actually provided.

Table 7 displays the total proportion of different anomalies' occurrences in the Testbed and presents the detection rate of each anomalous pattern by GMM and HMM. Here we consider the anomalous test data and the highest likelihood thresholds of both models (0.8 for GMM and -10 for HMM) that provide the maximal detection rate. Detection ratio is determined by the overall number of time-slots marked as anomaly by the model divided by the total number of time-slots encounter particular types of anomaly. Comparing GMM and HMM once more demonstrates the superior capability of HMM in recognition of anomalous events, while providing unnoticeable false positive rate (Table 6). Among the various types of anomalies generated for the Testbed, the highest detection rate belongs to heavy usage pattern, producing by multiple users and then single user. The lowest detection ratio, however, originates from jamming channel with low interval. Although there are some specific anomalous instances that

are never detected by the model, regardless of the cut-off threshold, the overall detection rate of the HMM is quite satisfactory. We intend to improve the detection estimate capacity of HMM in our future work by proposing more complex variations of HMMs.

7.4 Conclusions

In this section we described the Testbed deployment on a single AP and the process of data collection from a RADIUS server. We further explained the anomalies generated deliberately to prepare the ground truth data for model evaluation. We reproduced AP Shutdown/Halt, Heavy Usage from a single user and multiple users, and various types of interferences as a set of network anomalies. We then applied our proposed model to detect anomalous points, and discussed the effectiveness of each model. The experimental results demonstrated that HMM outperformed GMM in obtaining higher detection ratio while producing minor false alarm.

8 Conclusions and Future Work

In large deployments of 802.11 networks with varying usage, channel conditions, and operational constraints, network managers often demand tools that provide them with a comprehensive view of the entire network. Analyzing the users' behavioral patterns, learning efficient models to detect anomalous periods, and measuring the temporal performance of the network under certain circumstances are of great significance to provide an adequate level of satisfaction for the wireless users. Proposing time-invariant and time-variant modeling approaches and utilizing those models for anomaly detection in addition to a RADIUS testbed deployment with simulated anomalies compose the key contributions of this work.

We proposed a new application of HMMs in performance anomaly detection of 802.11 wireless networks and explored the necessity of temporal specifications of HMM rather than its simple time-independent counterpart model, GMM. We performed analysis and compared HMM and GMM in terms of modeling competency and anomaly detection performance on the large FEUP dataset as well as a similar but minor version of the deployed testbed with provoked anomalies for evaluation purposes.

The experimental results show that HMM models are capable of detecting a great portion of provoked anomalies on unobserved test data set (up to 75% TPR), and even disclosing unintentional anomalies occurred during the normal days of experiment. Besides, the false positive ratio is fairly low (only 1.1%) in HMM that outperforms GMM both in detection and fall-out rate.

In future work we intend to propose an anomaly detection algorithm that works in unsupervised mode regardless of the anomalous information provided for the data records. Furthermore, we will propose more complex HMMs to characterize and distinguish various anomaly-related patterns. We also plan to extend the testbed to multiple APs to explore new aspects of anomalies that concern the mobility effects of the wireless users in AP vicinities.

Acknowledgements This work is financed by the ERDF European Regional Development Fund through the Operational Programme for Competitiveness and Internationalisation - COMPETE 2020 Programme within project POCI-01-0145-FEDER-006961, and by National Funds through the FCT Fundao para a Cincia e a Tecnologia (Portuguese Foundation for Science and Technology) as part of project UID/EEA/50014/2013. The first author is also sponsored by FCT grant SFRH/BD/99714/2014.

References

1. The FreeRADIUS Project. <http://freeradius.org/>. Accessed in February 2016.
2. The Internet Engineering Task Force (IETF). <https://www.ietf.org/>. Accessed in January 2016.
3. Rfc 2865 radius authentication. <http://tools.ietf.org/html/rfc2865>. Accessed in January 2010.
4. Rfc 2866 radius authentication. <http://tools.ietf.org/html/rfc2866>. Accessed in January 2016.
5. System Sciences at Isis. <http://systems-sciences.uni-graz.at/etextbook/bigdata/confusionmatrix.html>. Accessed in April 2016.
6. Wifijammer. <https://github.com/DanMcInerney/wifijammer>. Accessed in February 2016.
7. Atul Adya, Paramvir Bahl, Ranveer Chandra, and Lili Qiu. Architecture and techniques for diagnosing faults in ieee 802.11 infrastructure networks. In *Proceedings of the 10th Annual International Conference on Mobile Computing and Networking, MobiCom '04*, pages 30–44, New York, NY, USA, 2004. ACM.
8. Ihsan Akbar, William H Tranter, et al. Dynamic spectrum allocation in cognitive radio using hidden markov models: Poisson distributed case. In *SoutheastCon, 2007. Proceedings. IEEE*, pages 196–201. IEEE, 2007.
9. Anisa Allahdadi, Ricardo Morla, Ana Aguiar, and Jaime S Cardoso. Predicting short 802.11 sessions from radius usage data. In *Local Computer Networks Workshops (LCN Workshops), 2013 IEEE 38th Conference on*, pages 1–8. IEEE, 2013.
10. Anisa Allahdadi, Ricardo Morla, and Jaime S Cardoso. Outlier detection in 802.11 wireless access points using hidden markov models. In *Wireless and Mobile Networking Conference (WMNC), 2014 7th IFIP*, pages 1–8. IEEE, 2014.
11. Wojciech Bednarczyk and Piotr Gajewski. Hidden markov models based channel status prediction for cognitive radio networks. *Session 4P6 RF and Wireless Communication*, page 2088, July 2015.
12. Yu-Chung Cheng, John Bellardo, Péter Benkő, Alex C. Snoeren, Geoffrey M. Voelker, and Stefan Savage. Jigsaw: Solving the puzzle of enterprise 802.11 analysis. In *Proceedings of the 2006 Conference on Applications, Technologies, Architectures, and Protocols for Computer Communications, SIGCOMM '06*, pages 39–50, New York, NY, USA, 2006. ACM.
13. D. Dujovne, T. Turletti, and F. Filali. A taxonomy of ieee 802.11 wireless parameters and open source measurement tools. *Communications Surveys Tutorials, IEEE*, 12(2):249–262, Second 2010.
14. Chris Fraley and Adrian E Raftery. Model-based clustering, discriminant analysis, and density estimation. *Journal of the American statistical Association*, 97(458):611–631, 2002.
15. Chris Fraley, Adrian E. Raftery, Thomas Brendan Murphy, and Luca Scrucca. *mclust Version 5.1 for R: Normal Mixture Modeling for Model-Based Clustering, Classification, and Density Estimation*, 2015.
16. Chittabrata Ghosh, Carlos Cordeiro, Dharma P Agrawal, and M Bhaskara Rao. Markov chain existence and hidden markov models in spectrum sensing. In *Pervasive Computing and Communications, 2009. PerCom 2009. IEEE International Conference on*, pages 1–6. IEEE, 2009.
17. M. Heusse, F. Rousseau, G. Berger-Sabbatel, and A. Duda. Performance anomaly of 802.11b. In *INFOCOM 2003. Twenty-Second Annual Joint Conference of*

the *IEEE Computer and Communications*. *IEEE Societies*, volume 2, pages 836–843 vol.2, March 2003.

18. Ankur Kamthe, Miguel A Carreira-Perpinán, and Alberto E Cerpa. M&m: multi-level markov model for wireless link simulations. In *Proceedings of the 7th ACM Conference on Embedded Networked Sensor Systems*, pages 57–70. ACM, 2009.
19. Syed A Khayam and Hayder Radha. Markov-based modeling of wireless local area networks. In *Proceedings of the 6th ACM international workshop on Modeling analysis and simulation of wireless and mobile systems*, pages 100–107. ACM, 2003.
20. Kaushik Lakshminarayanan, Srinivasan Seshan, and Peter Steenkiste. Understanding 802.11 performance in heterogeneous environments. In *Proceedings of the 2nd ACM SIGCOMM workshop on Home networks*, pages 43–48. ACM, 2011.
21. Ratul Mahajan, Maya Rodrig, David Wetherall, and John Zahorjan. Analyzing the mac-level behavior of wireless networks in the wild. In *ACM SIGCOMM Computer Communication Review*, volume 36, pages 75–86. ACM, 2006.
22. Dossa Massa and Ricardo Morla. Abrupt ending of 802.11 ap connections. In *Computers and Communications (ISCC), 2013 IEEE Symposium on*, pages 000348–000353. IEEE, 2013.
23. Dossa Massa and Ricardo Morla. Modeling 802.11 ap usage through daily keep-alive event counts. *Wireless networks*, 19(5):1005–1022, 2013.
24. Anthony J. Nicholson, Yatin Chawathe, Mike Y. Chen, Brian D. Noble, and David Wetherall. Improved access point selection. In *Proceedings of the 4th International Conference on Mobile Systems, Applications and Services, MobiSys '06*, pages 233–245, New York, NY, USA, 2006. ACM.
25. U. Paul, A. Kashyap, R. Maheshwari, and S.R. Das. Passive measurement of interference in wifi networks with application in misbehavior detection. *Mobile Computing, IEEE Transactions on*, 12(3):434–446, March 2013.
26. Pratap S Prasad and Prathima Agrawal. Movement prediction in wireless networks using mobility traces. In *Consumer Communications and Networking Conference (CCNC), 2010 7th IEEE*, pages 1–5. IEEE, 2010.
27. L. Rabiner and B.-H. Juang. An introduction to hidden markov models. *ASSP Magazine, IEEE*, 3(1):4–16, 1986.
28. Douglas Reynolds. Gaussian mixture models. *Encyclopedia of Biometrics*, pages 827–832, 2015.
29. A. Schliep, I. G. Costa, B. Georgi, C. Hafemeister, A. Schonhuth, and M. P. Mahmud. GHMM library. <http://ghmm.org>. Accessed in March 2016.
30. Anmol Sheth, Christian Doerr, Dirk Grunwald, Richard Han, and Douglas Sicker. Mojo: A distributed physical layer anomaly detection system for 802.11 wlans. In *Proceedings of the 4th international conference on Mobile systems, applications and services*, pages 191–204. ACM, 2006.
31. Vivek Shrivastava, Shravan K Rayanchu, Suman Banerjee, and Konstantina Papagiannaki. Pie in the sky: Online passive interference estimation for enterprise wlans. In *NSDI*, volume 11, pages 25–25, 2011.
32. Vamsi Krishna Tumuluru, Ping Wang, and Dusit Niyato. Channel status prediction for cognitive radio networks. *Wireless Communications and Mobile Computing*, 12(10):862–874, 2012.



Anisa Allahdadi received the B.Sc. in Computer Science from BIHE University (Bahá'í Institute for Higher Education), Iran in 2006 and M.Sc in Software Engineering from BIHE University, Iran in 2010. She is currently a researcher in the Center for Telecommunications and Multimedia at INESC TEC and pursuing her Ph.D in the MAP-i Doctoral Programme in the Faculty of Engineering of University of Porto. Her research interest include network management, probabilistic modeling and anomaly detection in IEEE 802.11 based wireless networks.



Ricardo Morla is an assistant professor of electrical and computer engineering at the University of Porto and principal investigator at INESC Porto, Portugal. His research interests are in the area of automatic system management with an emphasis on probabilistic modeling, prediction, anomaly detection, and root-cause analysis for ICT systems including network infrastructure and smart environments. He holds a Ph.D. in computer science from Lancaster University UK.