

Using data mining to study the impact of topology characteristics on the performance of wireless mesh networks

Tânia Calçada
INESC Porto,
Faculdade de Engenharia,
Universidade do Porto,
Porto, Portugal
tcalcada@inescporto.pt

Paulo Cortez
Centro Algoritmi,
Dep. Sistemas de Informação,
Universidade do Minho,
Guimarães, Portugal
pcortez@dsi.uminho.pt

Manuel Ricardo
INESC Porto,
Faculdade de Engenharia,
Universidade do Porto,
Porto, Portugal
mricardo@inescporto.pt

Abstract—This paper quantifies the impact of topological characteristics on the performance of single radio multichannel IEEE 802.11 mesh networks. Topological characteristics are the number of nodes per subnetwork, the hop count, the neighbor node density, the hidden nodes, the number of nodes in the neighborhood of the gateway, and the hidden nodes in the neighborhood of the gateway. Network performance metrics are throughput, fairness and delay. The data mining Support Vector Machine (SVM) model was used to extract the relationships between the network topology metrics and the network performance metrics based on data results obtained through ns-2 simulation of random networks. The results obtained can be used as a basis to design channel assignment algorithms or to aid the deployment and management of single radio wireless mesh networks.

Index Terms—Mesh networks, multi-channel, channel assignment, single-radio, topology metrics, node density, hidden nodes, miss ratio, data mining, SVM, relative importance

I. INTRODUCTION

Wireless mesh networks are an efficient and low cost solution to expand Internet coverage to areas where infrastructure connections to IEEE 802.11 access points are hard or expensive. Wireless networks such as these are expected to have high levels of usage, demanding solutions such as multi-channel communications between the mesh access points to improve the throughput of wireless mesh networks.

In our scenario, mesh nodes are expected to have two wireless cards with independent of-the-shelf radio interfaces running the standard MAC 802.11 protocol on different frequency bands; one radio interface operates as an access point, leaving the other to be used to interconnect the node to the mesh network. This configuration avoids the interference between the radio interfaces as described in [1], [2], [3]. The scenario studied in this paper addresses single-radio mesh networks. Special mesh nodes are wireline connected to network infrastructures acting as gateways to the Internet.

This work was co-supported by the SitMe project from QREN-ON.2 program and FCT SFRH/BD/13444/2003.

As mentioned in [4], [5], the multi-channel approach to be implemented in single-radio mesh networks assigns all nodes on a path to a common channel, which creates multiple subnetworks if multiple channels are used. The overall network performance depends on the subnetworks performance which in turn depends on topology characteristics. In [1], metrics related to topological characteristics are identified and an evaluation of their impact on the network performance is presented for a set of arbitrary mesh networks. The topology metrics identified were (1) the number of nodes per subnetwork, (2) the mean hop count, (3) the neighbor node density, (4) the hidden nodes measured by *miss ratio*, (5) the number of nodes in the neighborhood of the gateway, and (6) the hidden nodes in the neighborhood of the gateway measured by *miss ratio*. The analysis in [1] gave important guidelines for the relative importance of network topology characteristics in network performance; however, the scenarios studied were few and arbitrarily selected. This paper studies an extended set of random channel assignment schemes, aimed to bring a better understanding of the impact network topology characteristics have on network performance.

The methodology used in this work is presented in Fig. 1. A set of 3500 topologies with 36 randomly positioned nodes was created using the network simulator ns-2 [6]. Two channel assignment schemes were then applied to each of the 3500 topologies, assigning one of two possible channels to each node. Each of the 7000 networks were simulated four times using ns-2 with two possible traffic loads and two different simulation seeds. The results from the simulated networks were analyzed using a python script to retrieve performance metrics and the corresponding topology metrics. The metrics data from the 28000 network simulations were used to train a data mining model using Support Vector Machines (SVM) [7] using the rminer library of the R tool [8]. The input parameters of the models are the six topology metrics enumerated above and the output of each model is one of the three performance metrics studied in this paper: network aggregate throughput, fairness, and delay. Using fitted data mining models, the effect

of topology metrics on performance was quantified using a sensitivity analysis procedure as implemented in rminer.

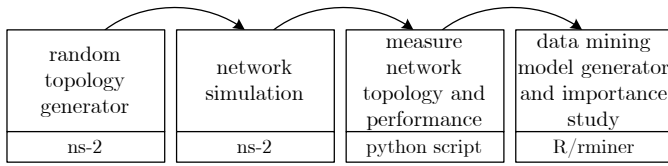


Fig. 1. Methodology.

This work provides 2 main contributions:

- 1) Quantification and ranking of the impact of topology metrics on the performance of wireless mesh networks, which are (by order of importance): (a) number of nodes in the 1st ring, i.e. nodes directly connected to a gateway; (b) mean hop count; (c) varying node quantity between subnetworks sharing a channel; (d) *miss ratio* in the 1st ring, which synthesizes the number of hidden nodes on the neighborhood of the gateway; (e) neighbor node density, which is the mean number of neighbors of a mesh node; (f) *miss ratio*, which synthesizes the number of hidden nodes on the network;
- 2) Use of data mining techniques for evaluation of the impact of wireless network topology characteristics on the performance of wireless mesh networks.

The rest of the paper is organized as follows: Section II describes the methodology used in this work, Section III identifies the performance metrics and the topology metrics that are used in the study, Section IV presents and explains the results obtained, and Section V concludes the paper.

II. METHODOLOGY

A. Random network generation

3500 different topologies were created using ns-2. Each network topology was generated randomly with consideration that the resultant network graph must be connected, signifying that every node in the network has multi-hop connectivity to all other nodes. Two nodes were randomly selected as gateways and operated on different radio channels. Each network has 36 nodes (including the gateways) and are spread randomly in a area of 1000 m × 1000 m.

Two different channel assignment schemes were applied to each topology. In the first assignment a random channel was assigned to each node guaranteeing multihop connection to a gateway. In the second channel assignment, the channel with a gateway closer in terms of hop count was assigned to each node; if both gateways were at the same distance, then the channel was selected randomly. Fig. 2 represents 2 instances of the generated networks and the lines between nodes represent wireless connectivity between them.

B. Network simulation

Each generated network was simulated four times using ns-2.29. Each node, not including the gateways, generated a traffic flow of either 480 kbit/s or 4.8 Mbit/s. These packets generate

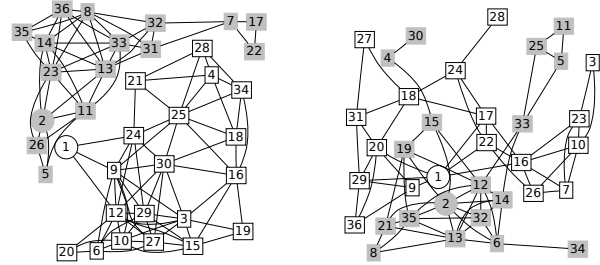


Fig. 2. Examples of generated random network topologies.

a Poisson distribution as they are User Datagram Protocol (UDP) and are destined to a node in the Internet through the serving gateway. Simultaneously, a node outside the mesh network generates a similar flow destined to each node in the network (gateways not included). All flows were configured with similar parameters which are fixed for each simulation each having two different seeds. The parameters used in each simulation are presented in Table I. The simulation tool ns-2 was used with a two-ray propagation model in the physical layer and MAC DCF 802.11 in the link layer. The Hybrid Wireless Mesh Protocol (HWMP) [9] was used to establish routes as defined in the IEEE standard to mesh networks.

Parameter	Value
Propagation Model	two ray ground
Channel data rate	54 Mbit/s
RX Threshold	-70.2 dBm, 350 m
Node distance	176 m
Packet size	1500 bytes
RTS/CTS	ON
Routing	HWMP
Source type	Poisson (UDP)
WarmUp	10 packet/s × 256 byte

TABLE I
PARAMETERS USED IN NS-2.29 SIMULATIONS.

The duration of each simulation was configured to give time to generate 10^4 packets on each flow with the exact duration depending on the rate of data flow. During the first 3 seconds there are no data flows to allow the HWMP routing protocol to execute the proactive tree building functionality. In this phase a route to one of the gateways is added to each node as described in the Proactive Path Request (PREQ) mechanism [9] and the reverse path is also created. Between 3 and 10 seconds the warm up flow takes place between each node and the gateway. Only the main flow packets are used to calculate the network performance.

C. Measure network topology and performance

To calculate the network topology and performance metrics, the two subnetworks resultant from the channel assignment are treated as a single network. The metrics aggregate the performance and the topology characteristics of all nodes in the network, independently of the channel the node is configured in.

D. Data mining model

The SVM was initially proposed by [7] for classification tasks (i.e., to model a discrete labeled output). After the introduction of the ε -insensitive loss function, it was possible to apply SVM to regression tasks [10]. SVM has theoretical advantages over other data mining techniques, such as the absence of local minima in the learning phase, i.e., the model always converges to the optimal solution. The main idea of the SVM is to transform the input data into a high-dimensional feature space by using a nonlinear mapping ϕ . Then, the SVM finds the best hyperplane within the feature space. The transformation depends on the kernel function adopted. The Gaussian kernel is the most popular one, presenting less parameters than other kernels, and thus is adopted in this work:

$$k(x, x') = e^{-\gamma \times \|x - x'\|^2}, \gamma > 0 \quad (1)$$

Under this setup, performance of the regression is affected by three parameters: γ – the parameter of the kernel, C – penalty parameter, and ε – the width of a ε -insensitive zone (both C and ε are used by SVM to select the support vectors during the learning phase). To reduce the search space, the first two values will be set using the heuristics of [11]: $C=3$ (for a standardized output) and $\varepsilon = \hat{\sigma}/\sqrt{N}$, where $\hat{\sigma} = 1.5/N \times \sum_{i=1}^N (y_i - \hat{y}_i)^2$ and \hat{y}_i is the value predicted by a 3-nearest neighbor algorithm. To optimize the most relevant SVM hyper-parameter (γ), we adopted a grid search under the range $\{2^{-15}, 2^{-13}, \dots, 2^3\}$, and an internal (i.e. over the training data) 3-fold cross validation was used to select the best γ value (i.e. that produces the lowest absolute deviation error on the validation data produced by the 3-fold scheme) [12]. After setting γ the SVM was retrained with all training data.

In order to evaluate the performance of the SVM predictions, we considered two popular regression metrics: Mean Absolute Error (MAD), and Coefficient of determination (R^2). Let y denote the target value, \hat{y} the predicted value, \bar{y} and $\bar{\hat{y}}$ the mean of these variables. Then:

$$\begin{aligned} MAD &= \sum_{i=1}^N |y_i - \hat{y}_i| / N \\ R^2 &= 1 - \left(\sum_{i=1}^N (y_i - \hat{y}_i)^2 / \sum_{i=1}^N (y_i - \bar{y})^2 \right) \end{aligned} \quad (2)$$

Lower values of MAD correspond to a higher predictive capacity, while the R^2 should be close to the unit value. To estimate the predictive performances of the SVM model, we applied 5 runs of an external (i.e. over all data) 5-fold cross validation. The predictive errors (i.e. MAD and R^2) shown in this work are reported in terms for the mean values of these runs and computed over the test (i.e. unseen) data defined by the 5-fold procedure. All experiments reported were implemented using the rminer library of the R tool [8].

We also apply the fitted SVM data mining models to estimate the impact of topology metrics on the wireless network performance. Despite the high complexity this model (due to the nonlinear transformation), it is still possible to extract knowledge in terms of input variable importance and Variable Effect Characteristic (VEC) curves by using a 1-D sensitivity

analysis [13]. This sensitivity analysis works by successively holding all inputs to their average values except one input, which is varied through its range of values in order to observe its effect on the target responses. The higher the variance observed in the responses, the higher is the relative importance of the input variable.

III. NETWORK METRICS

A. Network performance metrics

Network performance can be characterized using measures such as throughput, delay, packet loss, and number of re-transmissions. In our study we focus on node throughput and delay. First we want to maximize the average node throughput. Second, we want to maximize fairness among nodes' throughput to be sure that each mesh node offers an effective connection of its clients to the infrastructured network in order to avoid node starvation while assuring fair opportunities for packets transmission. Third, we want to minimize the end-to-end delay experienced by packets transmitted between the mesh nodes and the infrastructured network. We define the following metrics:

1) *Aggregate throughput*: The sum of the bit rate received by destinations. Formally, the aggregate throughput T_A measured on channel A is given in bits per second by Eq. 3 where T_i^{RX} is the number of packets received by node i , T_i^{TX} is the number of packets received by the gateway sent by node i , and L is the packet length given in bits.

$$T_A = \frac{\sum (T_i^{RX} + T_i^{TX})L}{\text{duration of simulation}} \quad (3)$$

The histogram of the aggregate throughput found in the simulations is presented in Fig. 3 for data rates of 480 kbit/s and 4.8 Mbit/s. When the flow data rate is low, the achieved aggregate throughput is also low. Experiences with higher flow data rate show a wider histogram meaning more uncertainty in the aggregate throughput results.

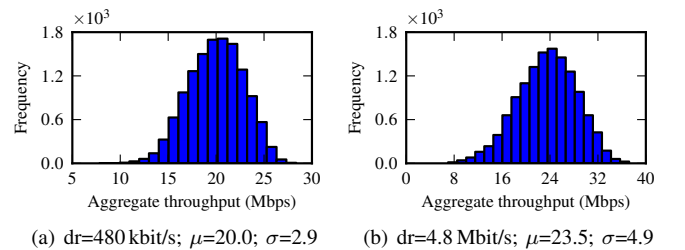


Fig. 3. Histogram for aggregate throughputs obtained by network simulations.

2) *Fairness*: The Jain Index [14] as given by Eq. 4 where T_i is the sum of T_i^{RX} and T_i^{TX} . The fairness J is independent of scale, applies to any number of nodes and is bounded between 0 and 1, where $J = 1$ indicate a totally fair network.

$$J = \frac{(\sum T_i)^2}{n \sum T_i^2} \quad (4)$$

The histogram of the fairness found in the simulations is presented in Fig. 4 for data rates of 480 kbit/s and 4.8 Mbit/s.

When the data rate is low, fairness is very high in most cases since all nodes in the mesh network have similar chances of transmitting their packets in an unsaturated network. When 34×2 flows (a download plus an upload flow per node) of 4.8 Mbit/s are placed on the network it inherently becomes saturated. Nodes that are not in the neighborhood of the gateways have more difficulty transmitting their packets through long multi-hop routes [1] and low values of fairness are achieved in most of the experiments.

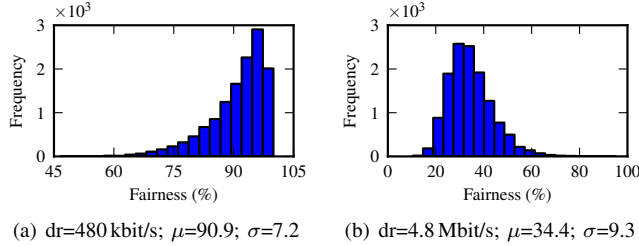


Fig. 4. Histogram for fairness obtained by the simulated networks.

3) *Delay*: The mean time elapsed between the creation of a packet and its reception at the final destination. Lost packets are not considered. The histogram of the delay measured in the simulations is presented in Fig. 5 for data rates of 480 kbit/s and 4.8 Mbit/s. When data rate is high (Fig. 5(b)), most packets to and from nodes that are not in the gateway neighborhood are lost [1]. Most of the packets considered to calculate the delay in high data rates scenarios (i.e. packets created close to the gateway) traversed less hops through the network to their final destination when compared to low data rate scenarios.

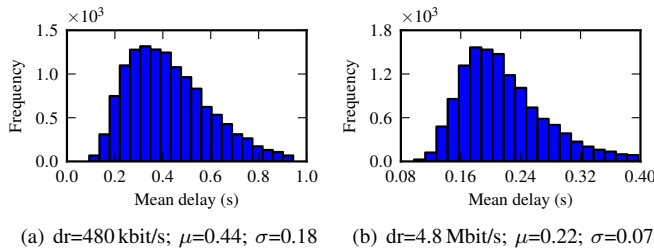


Fig. 5. Histogram for delay obtained by the simulated networks.

B. Topology metrics

The network topology metric considered to this study were previously presented in [1]. The number of nodes sharing a common channel is also considered since the network topology and channel assignment schemes are random.

1) *Number of nodes difference*: The difference between the number of nodes n in each subnetwork. Nodes in a subnetwork share a common radio channel and their packets are forwarded through a common gateway. For instance, if the subnetwork using channel A has 20 nodes and the subnetwork using channel B has 16 nodes, then the number of nodes difference is 4. The histogram of the number of nodes difference is presented in Fig. 6. Despite the scenarios being randomly generated, the number of nodes using either channel tend to be

equal since the channel and position of each node are selected randomly.

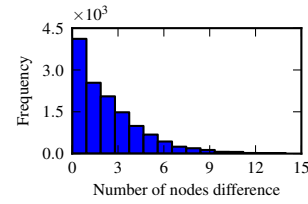


Fig. 6. Histogram for number of nodes difference between subnetworks.

2) *Neighbor node density*: The mean number of immediate neighbors each node has in the network (including the gateways). For example, the number of neighbors of node A is equal to the number of nodes in the receiving range (one hop away) of node A.

3) *Size of 1st ring*: The number of nodes a single hop distance from a gateway, which is also the neighbor node density that gateway.

4) *Mean hop count*: The mean value of the hop counts of each node in the network to reach a gateway (excluding the gateways and the nodes on the 1st rings which are included in the size of 1st ring metric).

5) *Miss ratio*: Defined in [15], is a global measure of the severity of hidden nodes in the overall network.

6) *1st ring Miss ratio*: defined in [1], it is a measure of the severity of the hidden nodes in the neighborhood of the gateway.

IV. RESULTS

The results obtained are presented in Fig. 7 and Fig. 8. The models obtained for low data rate scenarios have less errors (lower MAD values and higher R^2) than the models obtained with high loads, showing that non saturated networks maintain a more predictable behavior. Fig. 7 represents the relative importance of each topology metric for each model. For each model the total sum of all topology metrics relative importance is 1. Fig. 8 shows the VEC curves of each topology metrics for each model. The topology metrics were normalized to fit the [0,1] abscissa axis.

A. Model for aggregate throughput

Fig. 7(a) and Fig. 7(b) illustrates that the size of 1st ring is the main parameter for aggregate throughput models, with a relative importance of 0.52 and 0.72 for low and high traffic loads respectively. The VEC curves for aggregate throughput models presented in Fig. 8(a) and Fig. 8(b) show that high values of aggregate throughput are obtained when more nodes are in the neighborhood of the gateways (size of 1st ring). When a large number of nodes are directly connected to the gateway, there are two effects that contribute to increase the aggregate throughput: (1) data packets have to be forwarded only once in the mesh network until they reach the final destination, what means that less radio resources are used, leaving more opportunities for other packets to be transmitted;

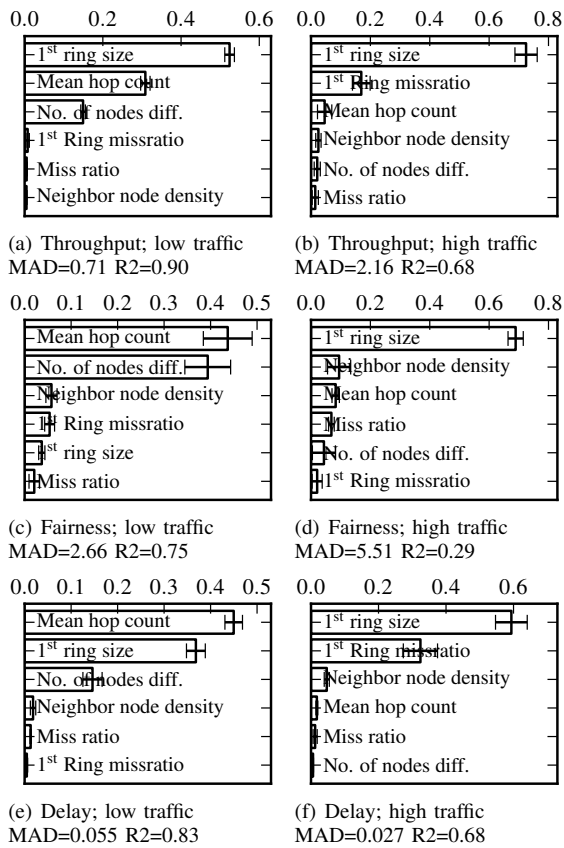


Fig. 7. Topology metrics importance on models.

(2) the gateway can effectively manage radio resources using the carrier sense and RTS/CTS, thus reducing the number of packet collisions.

When the network has low traffic loads, the metrics of mean hop count and the number of nodes difference metrics are key to the aggregate throughput model with relative importance of 0.30 and 0.15 respectively. Higher aggregate throughputs are obtained when nodes in the network are closer to the gateway (low mean hop count), as shown by Fig. 8(a). That result was expected as data packets have to be forwarded only once in the mesh network until they reach the final destination. Higher aggregate throughputs are obtained when the number of nodes in both subnetworks are similar, as shown by Fig. 8(a). When subnetworks have unbalanced number of nodes, the subnetwork with less nodes may obtain a high aggregate throughput, but the aggregate throughput of the overall network remains low. Neighbor node density, *Miss ratio* and *Miss ratio* on the 1st ring have relative importance values below 0.1 and are considered to have insignificant impact on the model output.

When the network has high traffic load, the mean hop count and the number of nodes difference do not significantly impact the aggregate throughput model, with relative importance of 0.04 and 0.02 respectively. This is mainly because only those nodes in the gateway neighborhood can transmit and receive packets as shown in [1] and also because of the low values of fairness shown in Fig. 4(b). This also explains why the

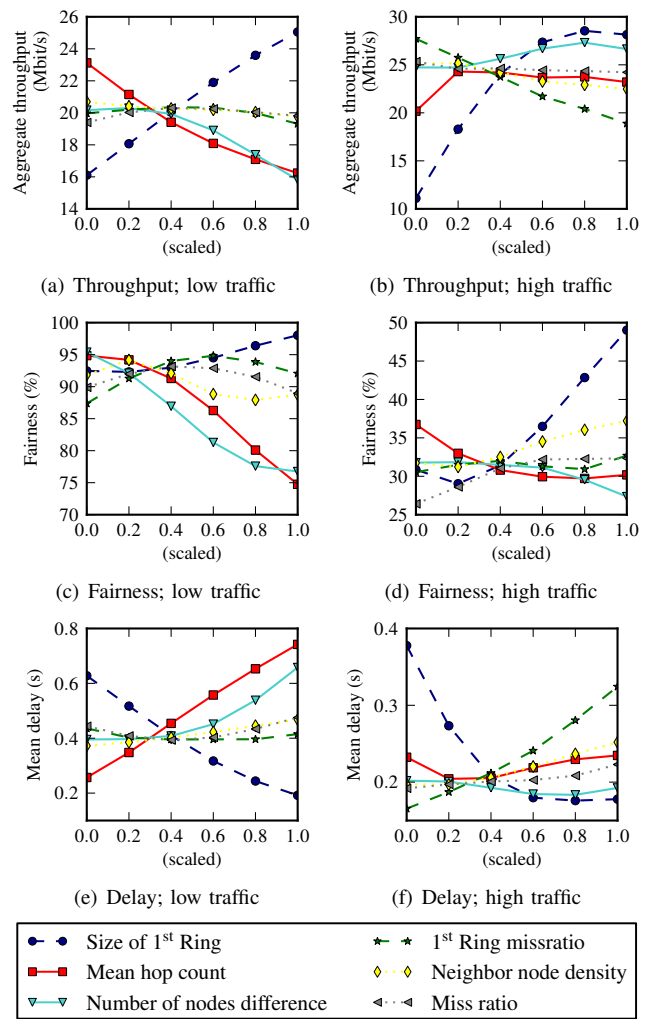


Fig. 8. VEC curves of models for high traffic load.

Miss ratio on the 1st ring has some impact in the aggregate throughput model when the network has high traffic load. High aggregate throughputs are obtained when links to the gateway have less hidden nodes. Neighbor node density and *Miss ratio* both have relative importance below 0.01 and are considered to have very little influence on the model output.

B. Model for fairness

Fairness models have high errors when compared with the throughput and delay models. This is because the topology metrics do not have a great effect on the network fairness as can be seen by VEC curves on Fig. 8(c) and Fig. 8(d). These curves show that even the topology metrics with higher relative importance values have a small variability on the output of the model. This data illustrates that fairness is an unpredictable performance metric in wireless mesh networks.

When the network has low traffic load, the mean hop count and the number of nodes difference are important to the fairness model with relative importance values of 0.44 and 0.39 respectively. On unsaturated mesh networks all nodes are able to transmit and receive their packets, therefore fairness is

close to 100% in most of the cases. However, if an unbalanced number of nodes exist in subnetworks using different channels, the overall network becomes unfair as shown in Fig. 8(c). The same applies when several nodes are more than 4 hops away from the gateway as shown in Fig. 8(c). Neighbor node density, *Miss ratio*, size of 1st ring, and *Miss ratio* on the 1st ring have relative importance values below 0.1 and are considered to have insignificant impact on the model output.

When the network has high traffic load, the size of 1st ring is the most influential parameter of the fairness model with a relative importance of 0.69. Under these traffic load conditions, the network becomes saturated and only those nodes in the neighborhood of the gateway are able to transmit and receive their packets successfully. When more nodes are close to the gateway the network becomes fairer as shown in Fig. 8(d).

C. Model for delay

Delay and aggregate throughput are related performance metrics as shown by the joint probability function plot on Fig. 9. Low aggregate throughputs occur when the delay is high, that is when packets take more time to reach their final destination, radio resources have less time available to transmit other packets resulting in lower aggregate throughputs. The models of these two performance metrics share all the important input parameters as can be observed by comparing figures 7(a), 7(b), 8(a), and 8(b) with 7(e), 7(f), 8(e), and 8(f). The main difference between these two models is that, for the delay model, under low traffic load conditions, the mean hop count have higher relative importance value than size of 1st ring. This result is expected, since for each hop that a packet performs, increases the overall transmission time due to the carrier sensing, which can eventually lead to packet retransmissions imposed by the MAC of IEEE 802.11 protocol.

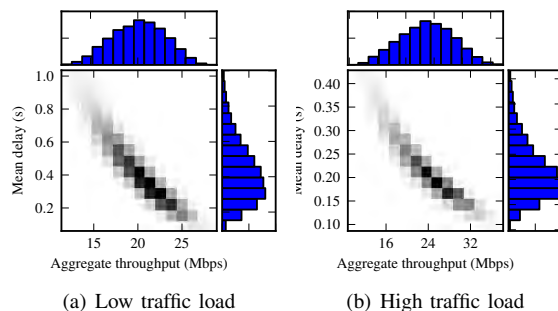


Fig. 9. Joint probability function plot of delay and aggregate throughput.

V. CONCLUSIONS

In this study we used a data mining approach based on the Support Vector Machine (SVM) algorithm and sensitivity analysis procedure to rank the relative importance of network topology characteristics on the performance of a single radio multichannel IEEE 802.11 mesh network. Experimental results suggest that the topology metric that has the greatest impact on the performance of a mesh network is the number of nodes that are directly connected to the gateway: a larger 1st ring size

resulted in increased data throughput, increased fairness and lower delay. Hop distance of nodes to the gateway and the varying number of nodes between subnetworks also impact network performance, but mostly in the realm of guaranteeing network fairness. As traffic load on the network approaches maximum capacity, network fairness begins to diminish for nodes far from the gateway while those nodes near the gateway experience enhanced levels of network fairness. The study also suggests that aggregate throughput could be improved by avoiding hidden nodes on links to the gateway.

SVMs have been successfully used to solve a wide range of complex problems in statistics, science and engineering. To the best of our knowledge, these techniques have not yet been applied to the study of wireless networks topologies or used to rank the impact of topology characteristics on wireless network performance.

As future work we plan to generalize our results by: (1) using network topologies with varying total number of nodes; (2) using v number varying channels; (3) using different traffic patterns and traffic loads; (4) modifying the output of data mining models by combining throughput, fairness and delay into a single utility function.

REFERENCES

- [1] T. Calçada and M. Ricardo, "The impact of network topology on the performance of multi-channel single-radio mesh networks," in *Proc. of Networking and Electronic Commerce Research Conf. (NAEC'11)*, (Italy), 2011.
- [2] F. Teixeira, T. Calçada, and M. Ricardo, "Protocol for channel assignment in single-radio mesh networks," in *Proc. Conf. on Mobile Networks and Management (MONAMI'11)*, (Portugal), 2011.
- [3] A. Dhananjay, H. Zhang, J. Li, and L. Subramanian, "Practical, distributed channel assignment and routing in dual-radio mesh networks," *ACM SIGCOMM Computer Communication Review*, vol. 39, no. 4, pp. 99–110, 2009.
- [4] J. So and N. Vaidya, "Load-balancing routing in multichannel hybrid wireless networks with single network interface," *IEEE Trans. Vehicular Technology*, vol. 55, no. 3, pp. 806–812, 2006.
- [5] R. Vedantham, S. Kakumanu, S. Lakshmanan, and R. Sivakumar, "Component based channel assignment in single radio, multi-channel ad hoc networks," in *Proc. Inter. Conf. on Mobile computing and networking (MobiCom'06)*, (USA), pp. 378–389, 2006.
- [6] "The network simulator ns-2." <http://www.isi.edu/nsnam/ns/>.
- [7] C. Cortes and V. Vapnik, "Support-vector networks," *Machine learning*, vol. 20, no. 3, pp. 273–297, 1995.
- [8] P. Cortez, "Data mining with neural networks and support vector machines using the r/rminer tool," *Advances in Data Mining. Applications and Theoretical Aspects*, pp. 572–583, 2010.
- [9] IEEE 802.11s, "IEEE wireless LAN medium access control (MAC) and physical layer (PHY) specifications draft on mesh networking," 2009.
- [10] A. J. Smola and B. Schölkopf, "A tutorial on support vector regression," *Statistics and computing*, vol. 14, no. 3, pp. 199–222, 2004.
- [11] V. Cherkassky and Y. Ma, "Practical selection of SVM parameters and noise estimation for SVM regression," *Neural Networks*, vol. 17, no. 1, pp. 113–126, 2004.
- [12] T. Hastie, R. Tibshirani, and J. H. Friedman, *The elements of statistical learning: data mining, inference, and prediction*. Springer Verlag, 2009.
- [13] P. Cortez and M. Embrechts, "Opening black box data mining models using sensitivity analysis," in *Proc. of IEEE Symp. on Computational Intelligence and Data Mining (CIDM'11)*, (France), 2011.
- [14] R. Jain, D. M. Chiu, and W. Hawe, "A quantitative measure of fairness and discrimination for resource allocation in shared systems," *Tech. Rep. 301*, Digital Equipment Corp., 1984.
- [15] L. B. Jiang and S. C. Liew, "Improving throughput and fairness by reducing exposed and hidden nodes in 802.11 networks," *IEEE Trans. Mobile Computing*, vol. 7, no. 1, pp. 34–49, 2008.