

Data Driven Hyperparameters Optimization of One-Class Support Vector Machines for Anomaly Detection in Wireless Sensor Networks

Van Vuong Trinh and Kim Phuc Tran

Division of Artificial Intelligence,

Faculty of Information Technology,

Dong A University, Danang, Vietnam

Email: vanvuong.trinh@gmail.com, phuctk@donga.edu.vn

Truong Thu Huong

Department of Telecommunication Systems,

School of Electronics and Telecommunications,

Hanoi University of Science and Technology, Hanoi, Vietnam

Email: huong.truongthu@hust.edu.vn

Abstract—One-class support vector machines (OCSVM) have been recently applied to detect anomalies in wireless sensor networks (WSNs). Typically, OCSVM is kernelized by radial basis functions (RBF, or Gaussian kernel) whereas selecting hyperparameters is based upon availability of labelled anomalous, which is rarely applicable in practice. This article investigates the application of OCSVM with data-driven hyperparameters optimization. Specifically, a kernel distance based optimization criteria is used instead of labelled data based metrics such as geometric mean accuracy (g-mean) or area under the receiver operating characteristic (AUROC). The efficiency of this method is illustrated over a real data set.

Index Terms—one-class support vector machines, anomaly detection, wireless sensor networks, Gaussian kernel, parameters selection.

I. INTRODUCTION

II. ONE-CLASS SUPPORT VECTOR MACHINES AND GAUSSIAN KERNEL

In this section, we briefly recall one-class support vector machines (OCSVM) [?]. OCSVM is used to estimate the support of a distribution. Notationally, let us consider a data set $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_i, \dots, \mathbf{x}_N\}$, with each $\mathbf{x}_i \in \mathcal{R}^D$ belonging to a given class of interest (named target class). The basic idea behind the OCSVM is to separate data from the origin by finding a hyperplane with maximum margin separation from the origin. In order to deal with nonlinearly problems, the hyperplane is defined in a high-dimensional Hilbert feature space \mathcal{F} where the samples are mapped through a nonlinear transformation $\Phi(\cdot)$. We will work only a kernel function $k(\mathbf{x}, \mathbf{y})$ instead of the scalar product $(\Phi(\mathbf{x}), \Phi(\mathbf{y}))$. To separate the data set from the origin, [1] solved the following quadratic program:

$$\underset{\mathbf{w} \in \mathcal{F}, \mathbf{a}, \boldsymbol{\xi} \in \mathcal{R}^N, \rho \in \mathcal{R}}{\text{Minimize}} \quad \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu N} \sum_{i=1}^N \xi_i - \rho \quad (1a)$$

$$\text{Subject to } (\mathbf{w} \cdot \Phi(\mathbf{x}_i)) \geq \rho - \xi_i, \quad \xi_i \geq 0 \quad \forall i = 1, \dots, N \quad (1b)$$

Here, \mathbf{w} is a vector perpendicular to the hyperplane in \mathcal{F} , and ρ is the distance to the origin. Since the training data

distribution may contain outliers, a set of slack variables $\xi_i \geq 0$ is introduced to deal with them. The parameter $\nu \in (0, 1]$ controls the tradeoff between the number of examples of the training set mapped as positive by the decision function

$$f(\mathbf{z}) = \text{sgn}((\mathbf{w} \cdot \Phi(\mathbf{x})) - \rho) \quad (2)$$

and having a small value of $\|\mathbf{w}\|$ to control model complexity.

Using multipliers $\alpha_i, \beta_i \geq 0$, [1] introduced a Lagrangian

$$L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \rho) = \frac{1}{2} \|\mathbf{w}\|^2 + \frac{1}{\nu N} \sum_{i=1}^N \xi_i - \rho - \sum_{i=1}^N \alpha_i ((\mathbf{w} \cdot \Phi(\mathbf{x}_i)) - \rho + \xi_i) - \sum_{i=1}^N \beta_i \xi_i \quad (3)$$

and set the derivatives with respect to the primal variables \mathbf{w} , $\boldsymbol{\xi}$, ρ equal to zero, i.e.

$$\frac{\partial L(\mathbf{w}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\beta}, \rho)}{\partial \mathbf{w}} = 0 : \quad \mathbf{w} = \sum_{i=1}^N \alpha_i \Phi(\mathbf{x}_i), \quad (4)$$

$$\frac{\partial L(R, \mathbf{a}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\gamma}, \rho)}{\partial \xi_i} = 0 : \quad \alpha_i = \frac{1}{\nu N} - \beta_i \leq \frac{1}{\nu N}, i = 1, \dots, N \quad (5)$$

$$\frac{\partial L(R, \mathbf{a}, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\gamma}, \rho)}{\partial \rho} = 0 : \quad \sum_{i=1}^N \alpha_i = 1 \quad (6)$$

In (4), all patterns $\{\mathbf{x}_i : i \in [1, \dots, N], \alpha_i > 0\}$ are called Support Vectors. From (4), using the kernel function, the decision function (2) is transformed into a kernel expansion

$$f(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^N \alpha_i k(\mathbf{x}_i, \mathbf{x}) - \rho \right) \quad (8)$$

Substituting (4), (5) and (6) into (3), and using the kernel function, we have

$$L(\boldsymbol{\alpha}) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j). \quad (9)$$

Then, we obtain the dual problem

$$\text{Minimize}_{\alpha} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j) \quad (10a)$$

$$\text{Subject to } 0 \leq \alpha_i \leq \frac{1}{\nu N}, \quad \forall i = 1 \dots N, \quad \sum_{i=1}^N \alpha_i = 1 \quad (10b)$$

[1] shown that at the optimum, the two inequality constraints (1b) became equalities if α_i and β_i are nonzero, which implies $0 < \alpha_i < \frac{1}{\nu N}$. Thus, the value of ρ can be recovered by exploiting that for any such α_i , the corresponding pattern \mathbf{x}_i satisfies

$$\rho = \langle \mathbf{w}, \Phi(\mathbf{x}_i) \rangle = \sum_{j=1}^N \alpha_j k(\mathbf{x}_j, \mathbf{x}_i) \quad (11)$$

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp \left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2} \right) \quad (12)$$

$$\begin{aligned} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 &= 2 - 2k(\mathbf{x}_i, \mathbf{x}_j) \\ &= 2 - 2 \exp \left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2} \right) \end{aligned} \quad (13)$$

III. DESCRIPTION OF ANOMALY DETECTION PROCEDURE FOR WSNs

[2]

$$\begin{aligned} J(\sigma) &= \frac{1}{n} \sum_{i=1}^N \max_j \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 \\ &\quad - \frac{1}{n} \sum_{i=1}^N \min_{j \neq i} \|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 \end{aligned} \quad (14)$$

$$J(\sigma) = \frac{2}{n} \sum_{i=1}^N \min_{j \neq i} k(\mathbf{x}_i, \mathbf{x}_j) - \frac{2}{n} \sum_{i=1}^N \max_j k(\mathbf{x}_i, \mathbf{x}_j) \quad (15)$$

$$\begin{aligned} J(\sigma) &= \frac{2}{n} \sum_{i=1}^N \exp \left(-\frac{\min_{j \neq i} \|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2} \right) \\ &\quad - \frac{2}{n} \sum_{i=1}^N \exp \left(-\frac{\max_j \|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2} \right) \end{aligned} \quad (16)$$

$$\begin{aligned} J(\sigma) &= \frac{2}{n} \sum_{i=1}^N \exp \left(-\frac{\text{Near}(\mathbf{x}_i)}{2\sigma^2} \right) \\ &\quad - \frac{2}{n} \sum_{i=1}^N \exp \left(-\frac{\text{Far}(\mathbf{x}_i)}{2\sigma^2} \right) \end{aligned} \quad (17)$$

$$\begin{aligned} \nabla J(\sigma) &= \frac{2}{n} \sum_{i=1}^N \exp \left(-\frac{\text{Near}(\mathbf{x}_i)}{2\sigma^2} \right) \frac{\text{Near}(\mathbf{x}_i)}{\sigma^3} \\ &\quad - \frac{2}{n} \sum_{i=1}^N \exp \left(-\frac{\text{Far}(\mathbf{x}_i)}{2\sigma^2} \right) \frac{\text{Far}(\mathbf{x}_i)}{\sigma^3} \end{aligned} \quad (18)$$

Gradient-based optimization

No convexity is guarantee

Time complexity of the DFN algorithm is $\mathcal{O}(N)$.

The whole procedure is summarized as below.

Algorithm 1 (Anomaly detection procedure) Assume a training set $\{\mathbf{x}_k^{\text{train}}\}_k$, a threshold δ and kernel $\kappa(\cdot, \cdot)$ are given. This algorithm produces the decision function $f(\cdot)$ defined by the constant D , the support vectors \mathbf{x}_i and corresponding Lagrange multipliers α_i .

▷ Training phase:

1 Set hyperparameters as

▷ Decision phase:

2 For a new sample \mathbf{z} , classify it according to (??), then raise an alarm if $f(\mathbf{z}) = 1$.

IV. ILLUSTRATIVE EXAMPLE

This section investigates the efficiency of anomaly detection algorithm over a real data set. The source code is freely available at <https://github.com/trinhvv/wsn-ocsvm-dfn>. All computation was performed on a platform with 2.6 GHz Intel(R) Core(TM) i7 and 16GB of RAM.

A. Data description

We consider a data set gathered from a WSN deployment at the Intel Berkeley Research Laboratory (IBRL) [3] with 54 *Mica2Dot* sensor nodes. Fig. 1 shows the sensor deployment in the laboratory. The sensors collect five measurements: light in Lux, temperature in degrees celsius, humidity (temperature corrected relative humidity) ranging from 0% to 100%, voltage in volts and network topology information in each 30 second sampling period. Node 0 is the gateway node while other nodes broadcast their data in multiple hops to the gateway node. During the 30 day (720 hour) period between 28th Feb 2004 and 5th April 2004, the 54 nodes collected about 2.3 million readings.

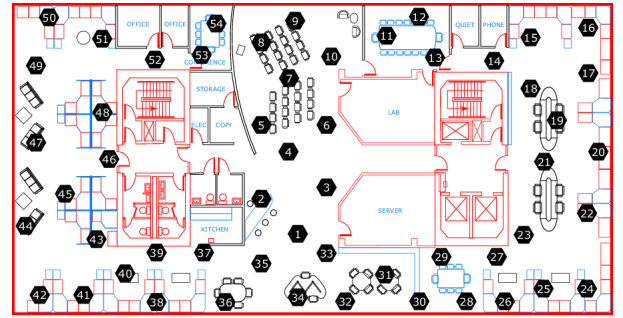


Fig. 1: A map of sensors' location. (Source: [3])

In this paper we consider the IBRL data set obtained from 5 close nodes, 1, 2, 33, 35, 37. Also, only two features, namely temperature and humidity, are taken into account. The data during the first 10 days period on March 2004 will be used as the training set. This training set contains more than 82000 samples.

In order to evaluate performance of the proposed method, we also use a testing set in some concrete time intervals. Since the original data did not contain any labels as to which data is normal and anomalous, we visually identify and label them as normal and anomalous. This data set contains about 10000 normal and 4000 anomalous samples.

B. Gaussian kernel parameter optimization

C. Training OCSVM and some results

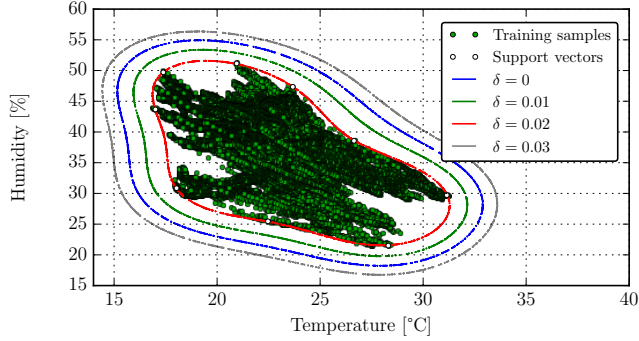


Fig. 2: Discrimination boundaries with some δ 's values.

δ	0.000	0.01	0.02	0.03
DR [%]	100	100	100	100
FPR [%]	0	0	0	0

TABLE I: DR and FPR versus δ .

V. CONCLUSION AND FUTURE WORK

REFERENCES

- [1] B. Schölkopf, J. C. Platt, J. Shawe-Taylor, A. J. Smola, and R. C. Williamson, "Estimating the support of a high-dimensional distribution," *Neural computation*, vol. 13, no. 7, pp. 1443–1471, 2001.
- [2] Y. Xiao, H. Wang, L. Zhang, and W. Xu, "Two methods of selecting Gaussian kernel parameters for one-class SVM and their application to fault detection," *Knowledge-Based Systems*, vol. 59, pp. 75–84, 2014. [Online]. Available: <http://dx.doi.org/10.1016/j.knosys.2014.01.020>
- [3] P. Buonadonna, D. Gay, J. M. Hellerstein, W. Hong, and S. Madden, "TASK: Sensor network in a box," *Proceedings of the Second European Workshop on Wireless Sensor Networks, EWSN 2005*, vol. 2005, pp. 133–144, 2005.

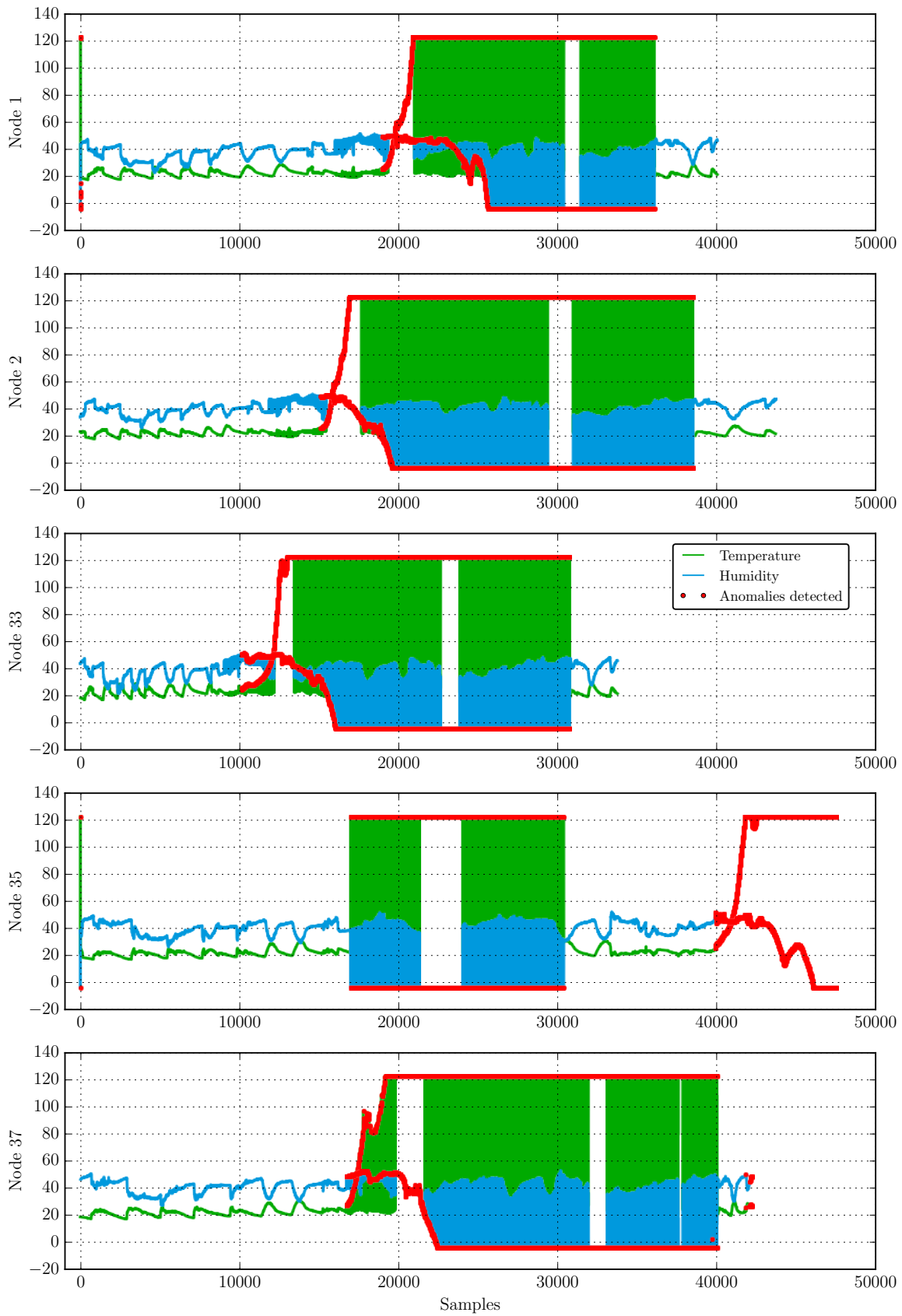


Fig. 3: Anomaly detection validation upon whole IBRL data set on 5 nodes. Almost apparent anomalies, i.e. temperature and humidity measurements that are too high or too low, are detected.