

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

Thu Huong Truong

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 PRELIMINARIES

In this section, we briefly recall one-class support vector machines (OCSVM) [1].

A. Theory

Given N samples \mathbf{x}_k , $k = 1, \dots, N$, SVDD method aims to estimate a sphere with minimum volume that contains all (or most of) these data. It is also assumed that almost these training samples belong to an unknown distribution. Let \mathbf{a} and R being reserved for the center and the radius of the sphere, we define the objective function to minimize the volume of the sphere and the number of outliers as:

$$R^2 + C \sum_{k=1}^N \xi_k \quad (1)$$

where $C > 0$ is a regularization parameter with constraints that almost data points are within the sphere:

$$\|\mathbf{x}_k - \mathbf{a}\|^2 \leq R^2 + \xi_k, \xi_k \geq 0 \quad \forall k \quad (2)$$

To adapt with nonspherical distribution, a conventional approach is to map given data into a higher dimensional feature

space, then learning a sphere in such a new space. This results into the so-called *primal optimisation* as follows:

$$\text{Minimize } R^2 + C \sum_{k=1}^N \xi_k \quad (3a)$$

$$\text{Subject to } \|\phi(\mathbf{x}_k) - \mathbf{a}\|^2 \leq R^2 + \xi_k, \xi_k \geq 0 \quad \forall k \quad (3b)$$

where $\phi(\cdot)$ is the aforementioned feature mapping. The Lagrangian is hereafter written as:

$$\begin{aligned} \mathcal{L} = & R^2 + C \sum_{k=1}^N \xi_k \\ & - \sum_{k=1}^N \alpha_k \left[R^2 + \xi_k - \|\phi(\mathbf{x}_k) - \mathbf{a}\|^2 \right] - \sum_{k=1}^N \gamma_k \xi_k \end{aligned} \quad (4)$$

with the Lagrange multipliers $\alpha_k, \gamma_k \geq 0$. \mathcal{L} should be minimized w.r.t. R, \mathbf{a}, ξ_k and maximized w.r.t. α_k, γ_k .

Setting partial derivatives w.r.t. R, \mathbf{a}, ξ gives:

$$\frac{\partial \mathcal{L}}{\partial R} = 0 : \quad \sum_{k=1}^N \alpha_k = 1 \quad (5a)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{a}} = 0 : \quad \mathbf{a} = \sum_{k=1}^N \alpha_k \phi(\mathbf{x}_k) \quad (5b)$$

$$\frac{\partial \mathcal{L}}{\partial \xi_k} = 0 : \quad \alpha_k + \gamma_k = C \quad \forall k \quad (5c)$$

Obviously, Lagrange multipliers γ_k can be eliminated by imposing bound constraints on α_k as:

$$0 \leq \alpha_k \leq C \quad \forall k \quad (6)$$

Substituting (5a)-(5c) into (4) leads to the following *dual*

optimisation:

$$\begin{aligned} \underset{\alpha}{\text{Maximize}} \quad & \sum_{k=1}^N \alpha_k (\phi(\mathbf{x}_k) \cdot \phi(\mathbf{x}_k)) \\ & - \sum_{k=1}^N \sum_{l=1}^N \alpha_k \alpha_l (\phi(\mathbf{x}_k) \cdot \phi(\mathbf{x}_l)) \end{aligned} \quad (7a)$$

$$\text{Subject to} \quad \sum_{k=1}^N \alpha_k = 1, 0 \leq \alpha_k \leq C \quad \forall k \quad (7b)$$

$$\|\phi(\mathbf{x}_k) - \mathbf{a}\|^2 < R^2 \rightarrow \alpha_k = 0 \quad (8a)$$

$$\|\phi(\mathbf{x}_k) - \mathbf{a}\|^2 = R^2 \rightarrow 0 < \alpha_k < C \quad (8b)$$

$$\|\phi(\mathbf{x}_k) - \mathbf{a}\|^2 > R^2 \rightarrow \alpha_k = C \quad (8c)$$

B. Kernelization

Instead of using inner product, an alternative kernel product can also be adopted:

$$\kappa(\mathbf{x}_k, \mathbf{x}_l) = \phi(\mathbf{x}_k) \cdot \phi(\mathbf{x}_l) \quad (9)$$

This is the known *kernel trick* [2], aims to avoid the need of explicitly declaring a feature mapping $\phi(\cdot)$.

$$\kappa(\mathbf{x}_k, \mathbf{x}_l) = \exp\left(-\frac{(\mathbf{x}_k - \mathbf{x}_l)'(\mathbf{x}_k - \mathbf{x}_l)}{2\sigma}\right) \quad (10)$$

where parameter σ is the kernel width.

$$\underset{\alpha}{\text{Maximize}} \quad \sum_{k=1}^N \alpha_k \kappa(\mathbf{x}_k, \mathbf{x}_k) - \sum_{k=1}^N \sum_{l=1}^N \alpha_k \alpha_l \kappa(\mathbf{x}_k, \mathbf{x}_l) \quad (11a)$$

$$\text{Subject to} \quad \sum_{k=1}^N \alpha_k = 1, 0 \leq \alpha_k \leq C \quad \forall k \quad (11b)$$

III. DESCRIPTION OF ANOMALY DETECTION PROCEDURE FOR WSNS

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.

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] B. Schölkopf, "The kernel trick for distances," *Advances in Neural Information Processing Systems 13*, vol. 13, pp. 301–307, 2001.