

GLAD: GLocalized Anomaly Detection via Active Feature Space Suppression

October 4, 2018

1 Introduction

Approach We denote the input space by $\mathcal{X} \subseteq \mathbb{R}^d$, and the complete set of available instances by \mathbf{D} . Each instance $\mathbf{x} \in \mathbf{D}$ is associated with a hidden label $y_i \in \{-1, +1\}$. Instances labeled $+1$ represent the *anomaly* class and are at most a small fraction τ of all instances. The label -1 represents the *nominal* class. We assume that we have an ensemble $\mathcal{E} = \{\mathcal{A}_1, \dots, \mathcal{A}_M\}$ of M pre-trained anomaly detectors. In order to customize an anomaly detector \mathcal{A}_m for a particular task, we start by decomposing the anomaly score it assigns to an instance $\mathbf{x} \in \mathcal{X}$ into two parts: **(a)** the original score $s_m(\mathbf{x})$, and **(b)** the *relevance* $p_m(\mathbf{x}) \in [0, 1]$. The score assigned by \mathcal{A}_m is then $s_m(\mathbf{x})p_m(\mathbf{x})$. The overall anomaly score is computed as: $\text{Score}(\mathbf{x}) = \sum_{m=1}^M s_m(\mathbf{x})p_m(\mathbf{x})$.

We start with the assumption that each ensemble member is uniformly relevant in every part of the input feature space, i.e., $p_1(\mathbf{x}) = \dots = p_M(\mathbf{x}) = \text{const} \ \forall \mathbf{x} \in \mathcal{X}$. This assumption is implemented by priming a neural network referred to as *FSSN* to predict the same probability value $b \in (0, 1) \ \forall \mathbf{x} \in \mathbf{D}$ using the ℓ_{prior} loss in Equation 1. In effect, this places ***a uniform prior for the relevance of each detector over the input space \mathcal{X}*** (rather than the *parameter* space). When all detectors have the same relevance, the final anomaly score simply corresponds to the average score across all detectors (up to a multiplicative constant). This is **a good starting point for active learning**. Next, the algorithm receives label feedback from an analyst and determines whether the ensemble made an error (i.e., ranked a labeled nominal in the top τ fraction of instances which are assumed to be anomalous, or ranked a labeled anomaly lower than the top τ fraction of instances). If so, the FSSN tries to suppress all erroneous detectors for similar inputs using a combination of ℓ_{prior} and ℓ_{AAD} (Equation 4).

$$\ell_{prior}(\mathbf{x}) = \sum_{m=1}^M -b \log(p_m(\mathbf{x})) - (1-b) \log(1-p_m(\mathbf{x})) \quad (1)$$

$$\ell_A(q; (\mathbf{x}, y)) = \max(0, y(q - \text{Score}(\mathbf{x}))) \quad (2)$$

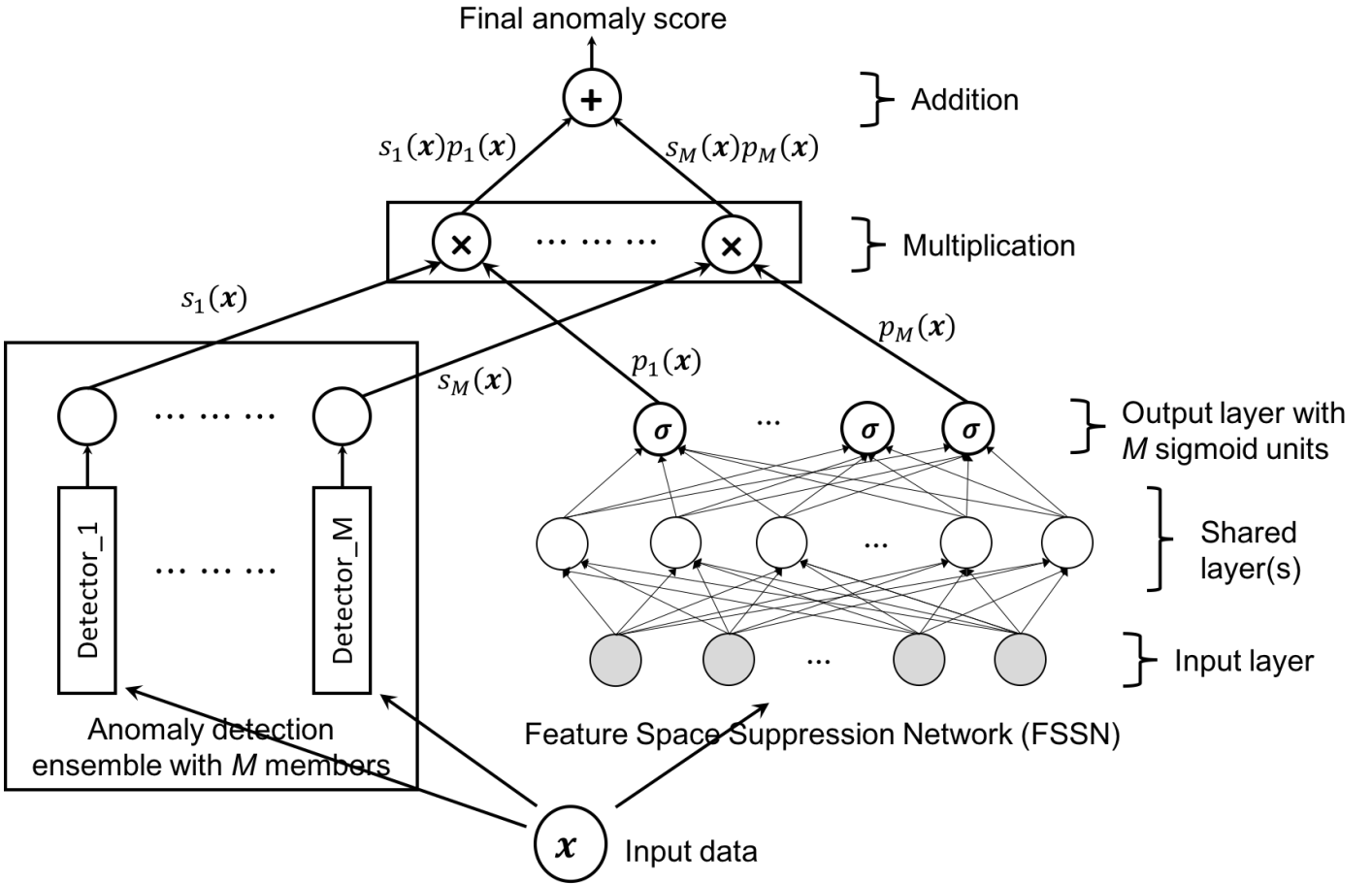
$$\ell_{AAD}(\mathbf{x}, y) = \ell_A(q_\tau^{(t-1)}; (\mathbf{x}, y)) + \ell_A(\text{Score}(\mathbf{x}_\tau^{(t-1)}); (\mathbf{x}, y)) \quad (3)$$

$$\ell_{FSSN} = \frac{1}{|\mathbf{H}_f|} \sum_{(\mathbf{x}, y) \in \mathbf{H}_f} \ell_{AAD}(\mathbf{x}, y) + \frac{\lambda}{|\mathbf{D}|} \sum_{\mathbf{x} \in \mathbf{D}} \ell_{prior}(\mathbf{x}) \quad (4)$$

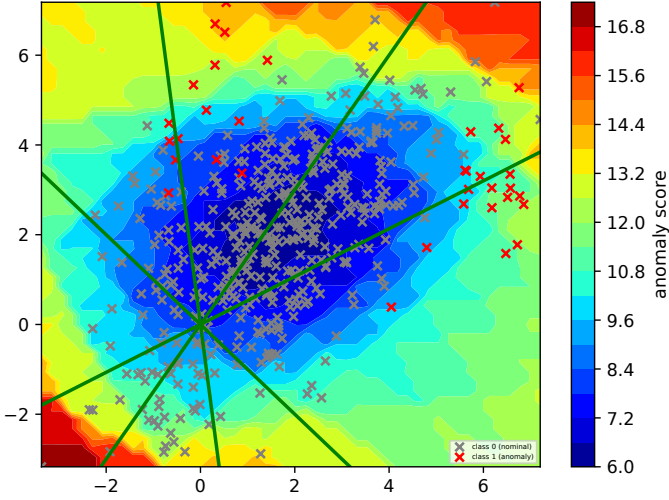
where $\lambda = 1$ works well

$\mathbf{H}_f \subseteq \mathbf{D}$ in Equation 4 denotes the set of instances labeled by the analyst. $\mathbf{x}_\tau^{(t-1)}$ and $q_\tau^{(t-1)}$ are the instance and the score resp. ranked at the τ -th quantile after the previous feedback iteration. ℓ_A encourages the scores of anomalies in \mathbf{H}_f to be higher than that of q , and the scores of nominals in \mathbf{H}_f to be lower.

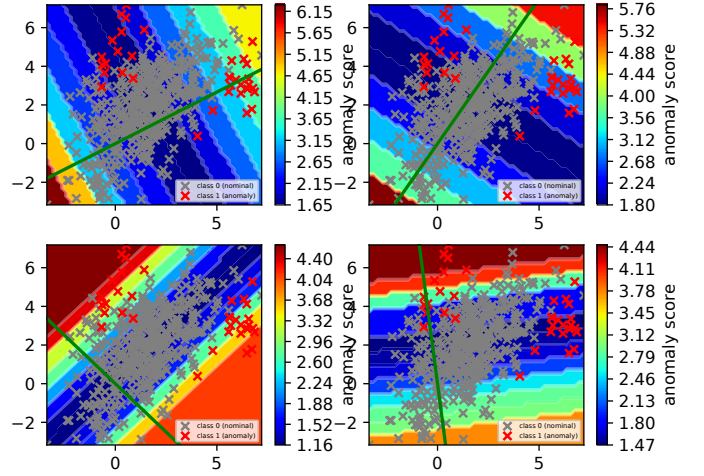
Toy Data The *Toy* dataset and the corresponding LODA ensembles have been shown below as an illustration of how GLAD works. We find that GLAD learns useful relevance information that can be of help to the analyst.



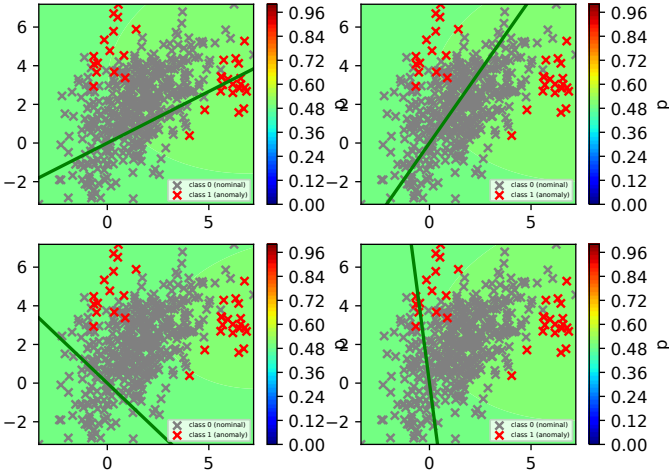
GLAD Architecture. The anomaly detection ensemble contains M detectors. We assume that these are all pre-trained and cannot be modified. The final layer of the Feature Space Suppression Network (FSSN) contains M sigmoid outputs, with each one paired with a corresponding ensemble member. Each output node in the FSSN is initially primed to predict the same probability (0.5 in our experiments) across the entire feature space. It then receives feedback from the analyst and learns which parts of the feature space are **relevant** for each detector. For an instance \mathbf{x} , $s_m(\mathbf{x})$ denotes the score assigned to it by the m -th detector, while $p_m(\mathbf{x})$ denotes the probability computed by the FSSN that the m -th detector is relevant. The final anomaly score for an instance \mathbf{x} is the sum of all its scores from each detector weighted by the corresponding relevances.



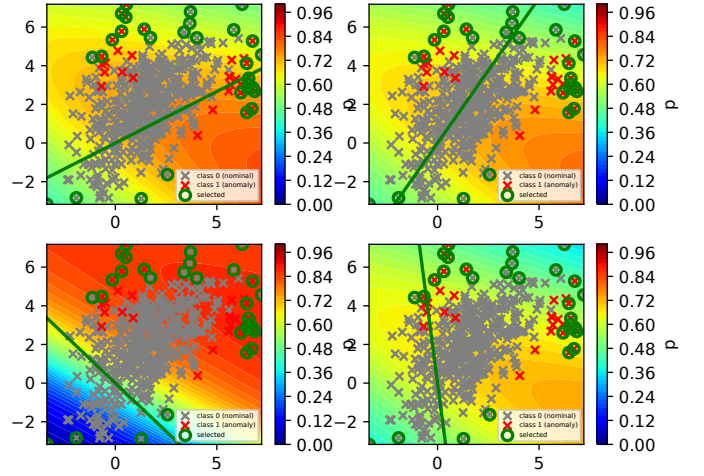
(a) Baseline LODA score contours



(b) Score contours for each LODA projection



(c) Initial uniform relevance ($b = 0.5$)



(d) Relevance after 30 feedback iterations

Toy data. More red on the **top row** indicates more *anomalous*. More red on the **bottom row** indicates more *relevant*. The red ‘x’ are true anomalies and grey ‘x’ are true nominals. (a) LODA with four projections (green lines) applied to the *Toy* dataset. (b) The contours of only the *bottom left* LODA projection are somewhat aligned with the true anomalies, i.e., most anomalies lie in the higher anomaly score regions. Other projections are highly inaccurate. (c) The output nodes of the FSSN are initially primed to return a relevance of 0.5 everywhere in the feature space. (d) The points circled in green were shown to the analyst for labeling, one per feedback iteration. After 30 iterations, the bottom left projection was found to be most relevant in the top-right half-space, whereas it is completely irrelevant in the bottom-left half-space. Other projections were less relevant in most parts of the feature space.