# ANOMALY DETECTION USING ONE-CLASS NEURAL NETWORKS

**Raghavendra Chalapathy**
University of Sydney,
Capital Markets Co-operative Research Centre (CMCRC)
rcha9612@uni.sydney.edu.au

**Aditya Krishna Menon**
Data61/CSIRO and the Australian National University
aditya.menon@data61.csiro.au

**Sanjay Chawla**
Qatar Computing Research Institute (QCRI), HBKU
schawla@qf.org.qa

January 14, 2019

## ABSTRACT

We propose a one-class neural network (OC-NN) model to detect anomalies in complex data sets. OC-NN combines the ability of deep networks to extract progressively rich representation of data with the one-class objective of creating a tight envelope around normal data. The OC-NN approach breaks new ground for the following crucial reason: data representation in the hidden layer is driven by the OC-NN objective and is thus customized for anomaly detection. This is a departure from other approaches which use a hybrid approach of learning deep features using an autoencoder and then feeding the features into a separate anomaly detection method like one-class SVM (OC-SVM). The hybrid OC-SVM approach is sub-optimal because it is unable to influence representational learning in the hidden layers. A comprehensive set of experiments demonstrate that on complex data sets (like CIFAR and GTSRB), OC-NN performs on par with state-of-the-art methods and outperformed conventional shallow methods in some scenarios.

***Keywords*** one class svm, anomalies detection, outlier detection, deep learning

## 1 Anomaly detection: motivation and challenges

A common need when analysing real-world datasets is determining which instances stand out as being dissimilar to all others. Such instances are known as *anomalies*, and the goal of *anomaly detection* (also known as *outlier detection*) is to determine all such instances in a data-driven fashion Chandola et al. [2007]. Anomalies can be caused by errors in the data but sometimes are indicative of a new, previously unknown, underlying process; in fact Hawkins Hawkins [1980] defines an outlier as an observation that *deviates so significantly from other observations as to arouse suspicion that it was generated by a different mechanism.* Unsupervised anomaly detection techniques uncover anomalies in an unlabeled test data, which plays a pivotal role in a variety of applications, such as, fraud detection, network intrusion detection and fault diagnosis. One-class Support Vector Machines (OC-SVM) Schölkopf and Smola [2002], Tax and Duin [2004] are widely used, effective unsupervised techniques to identify anomalies. However, performance of OC-SVM is sub-optimal on complex, high dimensional datasets Vapnik and Vapnik [1998], Vishwanathan et al. [2003], Bengio et al. [2007]. From recent literature, unsupervised anomaly detection using deep learning is proven to be very effective Zhou and Paffenroth [2017], Chalapathy et al. [2017]. Deep learning methods for anomaly detection can be broadly classified into model architecture using autoencoders Andrews et al. [2016a] and hybrid models Erfani et al. [2016]. Models involving autoencoders utilize magnitude of residual vector (i,e reconstruction error) for making anomaly assessments. While hybrid models mainly use autoencoder as feature extractor, wherein the hidden layer representations are used as input to traditional anomaly detection algorithms such as one-class SVM (OC-SVM). Following the success of transfer learning Pan and Yang [2010] to obtain rich representative features, hybrid models

have adopted pre-trained transfer learning models to obtain features as inputs to anomaly detection methods. Although using generic pre-trained networks for transfer learning representations is efficient, learning representations from scratch, on a moderately sized dataset, for a specific task of anomaly detection is shown to perform better Andrews et al. [2016b]. Since the hybrid models extract deep features using an autoencoder and then feed it to a separate anomaly detection method like OC-SVM, they fail to influence representational learning in the hidden layers. In this paper, we build on the theory to integrate a OC-SVM equivalent objective into the neural network architecture. The OC-NN combines the ability of deep networks to extract progressively rich representation of data alongwith the one-class objective, which obtains the hyperplane to separate all the normal data points from the origin. The OC-NN approach is novel for the following crucial reason: data representation ,is driven by the OC-NN objective and is thus customized for anomaly detection. We show that OC-NN can achieve comparable or better performance in some scenarios than existing shallow state-of-the art methods for complex datasets, while having reasonable training and testing time compared to the existing methods.

We summarize our main contributions as follows:

- We derive a new one class neural network (OC-NN) model for anomaly detection. OC-NN uses a one class SVM like loss function to drive the training of the neural network.

- We propose an alternating minimization algorithm for learning the parameters of the OC-NN model. We observe that the subproblem of the OC-NN objective is equivalent to a solving a quantile selection problem.

- We carry out extensive experiments which convincingly demonstrate that OC-NN outperforms other state-of-the-art deep learning approaches for anomaly detection on complex image and sequence data sets.

The rest of the paper is structured as follows. In Section 2 we provide a detailed survey of related and relevant work on anomaly detection. The main OC-NN model is developed in Section 3. The experiment setup, evaluation metrics and model configurations are described in Section 4. The results and analysis of the experiments are the focus of Section 5. We conclude in Section 6 with a summary and directions for future work.

## 2 Background and related work on anomaly detection

Anomaly detection is a well-studied topic in Data Science Chandola et al. [2007], Aggarwal [2016]. Unsupervised anomaly detection aims at discovering rules to separate normal and anomalous data in the absence of labels. One-Class SVM (OC-SVM) is a popular unsupervised approach to detect anomalies, which constructs a smooth boundary around the majority of probability mass of data Schölkopf et al. [2001]. OC-SVM will be described in detail in Section 2.2. In recent times, several approaches of feature selection and feature extraction methods have been proposed for complex, high-dimensional data for use with OC-SVM Cao et al. [2003], Neumann et al. [2005]. Following the unprecedented success of using deep autoencoder networks, as feature extractors, in tasks as diverse as visual, speech anomaly detection Chong and Tay [2017], Marchi et al. [2017], several hybrid models that combine feature extraction using deep learning and OC-SVM have appeared Sohaib et al. [2017], Erfani et al. [2016]. The benefits of leveraging pre-trained transfer learning representations for anomaly detection in hybrid models was made evident by the results obtained, using two publicly available [1] pre-trained CNN models: ImageNet-MatConvNet-VGG-F (VGG-F) and ImageNet-MatConvNet-VGG-M (VGG-M) Andrews et al. [2016b]. However, these hybrid OC-SVM approaches are decoupled in the sense that the feature learning is task agnostic and not customized for anomaly detecion. Recently a deep model which trains a neural network by minimizing the volume of a hypersphere that encloses the network representations of the data is proposed Ruff et al. [2018], our approach differs from this approach by combining the ability of deep networks to extract progressively rich representation of data alongwith the one-class objective, which obtains the hyperplane to separate all the normal data points from the origin.

### 2.1 Robust Deep Autoencoders for anomaly detection

Besides the hybrid approaches which use OC-SVM with deep learning features another approach for anomaly detection is to use deep autoencoders. Inspired by RPCA Xu et al. [2010], unsupervised anomaly detection techniques such as robust deep autoencoders can be used to separate normal from anomalous data Zhou and Paffenroth [2017], Chalapathy et al. [2017]. Robust Deep Autoencoder (RDA) or Robust Deep Convolutional Autoencoder (RCAE) decompose input data $X$ into two parts $X = L_D + S$, where $L_D$ represents the latent representation the hidden layer of

---

[1]Pretrained-models:http://www.vlfeat.org/matconvnet/pretrained/.

the autoencoder. The matrix $S$ captures noise and outliers which are hard to reconstruct as shown in Equation 1. The decomposition is carried out by optimizing the objective function shown in Equation 1.

$$\min_{\theta, S} +||L_D - D_\theta(E_\theta(L_D))||_2 + \lambda \cdot \|S^T\|_{2,1} \tag{1}$$

$$s.t. \ \ X - L_D - S = 0$$

The above optimization problem is solved using a combination of backpropagation and Alternating Direction Method of Multipliers (ADMM) approach Boyd and Vandenberghe [2004]. In our experiments we have carried out a detailed comparision between OC-NN and approaches based on robust autoencoders.

## 2.2 One-Class SVM for anomaly detection

One-Class SVM (OC-SVM) is a widely used approach to discover anomalies in an unsupervised fashion Schölkopf and Smola [2002]. OC-SVMs are a special case of support vector machine, which learns a hyperplane to separate all the data points from the origin in a reproducing kernel Hilbert space (RKHS) and maximises the distance from this hyperplane to the origin. Intuitively in OC-SVM all the data points are considered as positively labeled instances and the origin as the only negative labeled instance. More specifically, given a training data $\mathbf{X}$, a set without any class information, and $\Phi(\mathbf{X})$ a RKHS map function from the input space to the feature space $F$, a hyper-plane or linear decision function $f(\mathbf{X}_{n:})$ in the feature space $F$ is constructed as $f(\mathbf{X}_{n:}) = w^T \Phi(\mathbf{X}_{n:}) - r$, to separate as many as possible of the mapped vectors $\Phi(\mathbf{X}_{n:}), n : 1, 2, ..., N$ from the origin. Here $w$ is the norm perpendicular to the hyper-plane and $r$ is the bias of the hyper-plane. In order to obtain $w$ and $r$, we need to solve the following optimization problem,

$$\min_{w, r} \frac{1}{2} \|w\|_2^2 + \frac{1}{\nu} \cdot \frac{1}{N} \sum_{n=1}^{N} \max(0, r - \langle w, \Phi(\mathbf{X}_{n:}) \rangle) - r. \tag{2}$$

where $\nu \in (0, 1)$, is a parameter that controls a trade off between maximizing the distance of the hyper-plane from the origin and the number of data points that are allowed to cross the hyper-plane (the false positives).

## 3 From One Class SVM to One Class Neural Networks

We now present our one-class Neural Network (OC-NN) model for unsupervised anomaly detection. The method can be seen as designing a neural architecture using an OC-SVM equivalent loss function. Using OC-NN we will be able to exploit and refine features obtained from unsupervised tranfer learning specifically for anomaly detection. This in turn will it make it possible to discern anomalies in complex data sets where the decision boundary between normal and anomalous is highly nonlinear.

### 3.1 One-Class Neural Networks (OC-NN)

We design a simple feed forward network with one hiden layer having linear or sigmoid activation $g(\cdot)$ and one output node. Generalizations to deeper architectures is straightforward. The OC-NN objective can be formulated as:

$$\min_{w, V, r} \frac{1}{2} \|w\|_2^2 + \frac{1}{2} \|V\|_F^2 + \frac{1}{\nu} \cdot \frac{1}{N} \sum_{n=1}^{N} \max(0, r - \langle w, g(V\mathbf{X}_{n:}) \rangle) - r \tag{3}$$

where $w$ is the scalar output obtained from the hidden to output layer, $V$ is the weight matrix from input to hidden units. Thus the key insight of the paper is to replace the dot product $\langle \mathbf{w}, \mathbf{\Phi}(\mathbf{X_{n:}}) \rangle$ in OC-SVM with the dot product $\langle \mathbf{w}, \mathbf{g}(\mathbf{VX_{n:}}) \rangle$. This change will make it possible to leverage transfer learning features obtained using an autoencoder and create an additional layer to refine the features for anomaly detection. However, the price for the change is that the objective becomes non-convex and thus the resulting algorithm for inferring the parameters of the model will not lead to a global optima.

### 3.2 Training the model

We can optimize Equation 3 using an alternate minimization approach: We first fix $r$ and optimize for $w$ and $V$. We then use the new values of $w$ and $V$ to optimize $r$. However, as we will show, the optimal value of $r$ is just the

$\upsilon$-quantile of the array $\langle w, g(V x_n)\rangle$. We first define the objective to solve for $w$ and $V$ as

$$\underset{w,V}{\operatorname{argmin}} \frac{1}{2}\|w\|_2^2 + \frac{1}{2}\|V\|_F^2 + \frac{1}{\upsilon} \cdot \frac{1}{N} \sum_{n=1}^{N} \ell(y_n, \hat{y}_n(w, V)) \tag{4}$$

where

$$\ell(y, \hat{y}) = \max(0, y - \hat{y})$$
$$y_n = r$$
$$\hat{y}_n(w, V) = \langle w, g(V x_n)\rangle$$

Similarly the optimization problem for $r$ is

$$\underset{r}{\operatorname{argmin}} \left( \frac{1}{N\upsilon} \cdot \sum_{n=1}^{N} \max(0, r - \hat{y}_n) \right) - r \tag{5}$$

**Theorem 1.** *Given $w$ and $V$ obtained from solving Equation 4, the solution to Equation 5 is given by the $\upsilon^{th}$ quantile of $\{\hat{y}_n\}_{n=1}^{N}$, where*
$$\hat{y}_n = \langle w, g(V x_n)\rangle.$$

*Proof.* We can rewrite Equation 5 as:

$$\underset{r}{\operatorname{argmin}} \left( \frac{1}{N\upsilon} \cdot \sum_{n=1}^{N} \max(0, r - \hat{y}_n) \right) - \left( r - \frac{1}{N} \sum_{n=1}^{N} \hat{y}_n \right)$$

$$= \underset{r}{\operatorname{argmin}} \left( \frac{1}{N\upsilon} \cdot \sum_{n=1}^{N} \max(0, r - \hat{y}_n) \right) - \left( \frac{1}{N} \sum_{n=1}^{N} [\![r - \hat{y}_n]\!] \right)$$

$$= \underset{r}{\operatorname{argmin}} \left( \sum_{n=1}^{N} \max(0, r - \hat{y}_n) \right) - \upsilon \cdot \left( \sum_{n=1}^{N} [\![r - \hat{y}_n]\!] \right)$$

$$= \underset{r}{\operatorname{argmin}} \sum_{n=1}^{N} [\![\max(0, r - \hat{y}_n) - \upsilon \cdot (r - \hat{y}_n)]\!]$$

$$= \underset{r}{\operatorname{argmin}} \sum_{n=1}^{N} \begin{cases} (1 - \upsilon) \cdot (r - \hat{y}_n) & \text{if } r - \hat{y}_n > 0 \\ -\upsilon \cdot (r - \hat{y}_n) & \text{otherwise} \end{cases}$$

We can observe that the derivative with respect to $r$ is

$$F'(r) = \sum_{n=1}^{N} \begin{cases} (1 - \upsilon) & \text{if } r - \hat{y}_n > 0 \\ -\upsilon & \text{otherwise.} \end{cases}$$

Thus, by F'( r ) = 0 we obtain

$$(1 - \upsilon) \cdot \sum_{n=1}^{N} [\![r - \hat{y}_n > 0]\!] = \upsilon \cdot \sum_{n=1}^{N} [\![r - \hat{y}_n \le 0]\!]$$

$$= \upsilon \cdot \sum_{n=1}^{N} (1 - [\![r - \hat{y}_n > 0]\!])$$

$$= \upsilon \cdot N - \upsilon \cdot \sum_{n=1}^{N} [\![r - \hat{y}_n > 0]\!],$$

or

$$\frac{1}{N} \sum_{n=1}^{N} [\![r - \hat{y}_n > 0]\!] = \frac{1}{N} \sum_{n=1}^{N} [\![\hat{y}_n < r]\!] = \upsilon \tag{6}$$

This means we would require the $\upsilon^{\text{th}}$ quantile of $\{\hat{y}_n\}_{n=1}^{N}$. $\qquad\square$

### 3.3 OC-NN Algorithm

We summarize the solution in Algorithm 1. We initialize $r^{(0)}$ in Line 2. We learn the parameters$(w, V)$ of the neural network using the standard Backpropogation(BP) algorithm (Line 7). In the experiment section, we will train the model using features extracted from an autoencoder instead of raw data points. However this has no impact on the OC-NN algorithm. As show in Theorem 3.1, we solve for $r$ using the $\upsilon$-quantile of the scores $\langle y_n \rangle$. Once the convergence criterion is satisfied, the data points are labeled normal or anomalous using the decision function $S_n = sgn(\hat{y}_n - r)$.

---

**Algorithm 1** one-class neural network (OC-NN) algorithm

---

1:  **Input:** Set of points $\mathbf{X}_{n:}, \quad n : 1, ..., N$
2:  **Output:** A Set of decision scores $S_{n:} = \hat{y}_{n:}$, n: 1,...,N for X
3:  Initialise $r^{(0)}$
4:  $t \leftarrow 0$
5:  **while** (no convergence achieved) **do**
6:      Find $(w^{(t+1)}, V^{(t+1)})$                                    $\triangleright$ Optimize Equation 4 using BP.
7:      $r^{t+1} \leftarrow \nu^{\text{th}}$ quantile of $\{\hat{y}_n^{t+1}\}_{n=1}^N$
8:      $t \leftarrow t + 1$
9:  **end**
10: Compute decision score $S_{n:} = \hat{y}_n - r$ for each $\mathbf{X}_{n:}$
11: **if** $(S_{n:} \geq 0)$ **then**
12:     $\mathbf{X}_{n:}$ is normal point
13: **else**
14:     $\mathbf{X}_{n:}$ is anomalous
15: **return** $\{S_n\}$

---

**Example:** We give a small example to illustrate that the minimum of the function

$$f(r) = \left( \frac{1}{N\upsilon} \cdot \sum_{n=1}^{N} \max(0, r - y_n) \right) - r$$

occurs at the the $\upsilon$-quantile of the set $\{y_n\}$.

Let $y = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$ and $\upsilon = 0.33$. Then the minimum will occur at $f(3)$ as detailed in the table below.

| $r$ | $f(r)(\text{expr})$ | $f(r)(\text{value})$ |
|---|---|---|
| 1 | $\frac{1}{9*.33}[0 + 0 + \ldots 0] - 1$ | $-1.00$ |
| 2 | $\frac{1}{9*.33}[1 + 0 + \ldots 0] - 2$ | $-1.67$ |
| 3 | $\frac{1}{9*.33}[2 + 1 + \ldots 0] - 3$ | $-1.99$ |
| 4 | $\frac{1}{9*.33}[3 + 2 + 1 + \ldots 0] - 4$ | $-1.98$ |
| 5 | $\frac{1}{9*.33}[4 + 3 + 2 + 1 \ldots 0] - 5$ | $-1.63$ |
| 6 | $\frac{1}{9*.33}[5 + 4 + 3 + 2 + 1 \ldots 0] - 6$ | $-0.94$ |
| 7 | $\frac{1}{9*.33}[6 + 5 + 4 + 3 + 2 + 1 \ldots 0] - 7$ | $0.07$ |
| 8 | $\frac{1}{9*.33}[7 + 6 + 5 + 4 + 3 + 2 + 1 \ldots 0] - 8$ | $1.43$ |
| 9 | $\frac{1}{9*.33}[8 + 7 + 6 + 4 + \ldots 0] - 9$ | $3.12$ |

## 4 Experimental Setup

In this section, we show the empirical effectiveness of OC-NN formulation over the state-of-the-art methods on real-world data. Although our method is applicable in any context where autoencoders may be used for feature representation, e.g., speech. Our primary focus will be on non-trivial high dimensional images.

### 4.1 Methods compared

We compare our proposed one-class neural networks (OC-NN) with the following state-of-the-art methods for anomaly detection:

| Dataset | # instances | # anomalies | # features |
|---------|-------------|-------------|------------|
| Synthetic | 190 | 10 | 512 |
| MNIST | single class | 1% ( from all class) | 784 |
| CIFAR−10 | single class | 10% ( from all class) | 3072 |
| GTSRB | 1050 (stop signs ) | 100 (boundary attack) | 3072 |

Table 1: Summary of datasets used in experiments.

- **OC-SVM -SVDD** as per formulation in Schölkopf and Smola [2002]

- **Isolation Forest** as per formulation in Liu et al. [2008].

- **Kernel Density Estimation (KDE)** as per formulation in Parzen [1962].

- **Deep Convolutional Autoencoder (DCAE)** as per formulation in Masci et al. [2011].

- **AnoGAN** as per formulation in Radford et al. [2015].

- **Soft-Bound and One Class Deep SVDD** as per formulation in Ruff et al. [2018].

- **Robust Convolutional Autoencoder (RCAE)** as per formulation in Chalapathy et al. [2017].

- **One-class neural networks (OC-NN)** [2], our proposed model as per Equation 3.

We used Keras Chollet et al. [2015] and TensorFlow Abadi et al. [2016] for the implementation of OC-NN, DCAE, RCAE [3] For OC-SVM[4] and Isolation Forest[5], we used publicly available implementations.

## 4.2 Datasets

We compare all methods on synthetic and four real-world datasets as summarized below :

- Synthetic Data, consisting of 190 normal data points 10 anomalous points drawn from normal distribution with dimension 512.

- MNIST, consisting of 60000 $28 \times 28$ grayscale images of handwritten digits in 10 classes, with 10000 test images LeCun et al. [2010].

- GTSRB , are $32 \times 32$ colour images comprising of adversarial boundary attack on stop signs boards Stallkamp et al. [2011].

- CIFAR−10 consisting of 60000 $32 \times 32$ colour images in 10 classes, with 6000 images per class Krizhevsky and Hinton [2009].

For each dataset, we perform further processing to create a well-posed anomaly detection task, as described in the next section.

## 4.3 Evaluation of Models

### 4.3.1 Baseline Model Parameters:

The proposed OC-NN method is compared with several state-of-the-art baseline models as illustrated in Table 3. The model parameters of shallow baseline methods are used as per implementation in Ruff et al. [2018]. Shallow Baselines (i) Kernel OC-SVM/SVDD with Gaussian kernel. We select the inverse length scale $\gamma$ from $\gamma \in 2^{-10}$ , $2^{-9}$, . . . , $2^{-1}$ via grid search using the performance on a small holdout set (10) % of randomly drawn test samples). We run all experiments for $\nu = 0.1$ and report the better result. (ii) Kernel density estimation (KDE). We select the bandwidth $h$ of the Gaussian kernel from $h \in 2^{0.5}, 2^1, ..., 2^5$ via 5-fold cross-validation using the log-likelihood score. (iii) For the Isolation Forest (IF) we set the number of trees to $t = 100$ and the sub-sampling size to $\Psi = 256$, as recommended in the original work Ruff et al. [2018].

[2]https://github.com/raghavchalapathy/oc-nn
[3]https://github.com/raghavchalapathy/rcae
[4]http://scikit-learn.org/stable/auto_examples/svm/plot_oneclass.html
[5]http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.IsolationForest.html
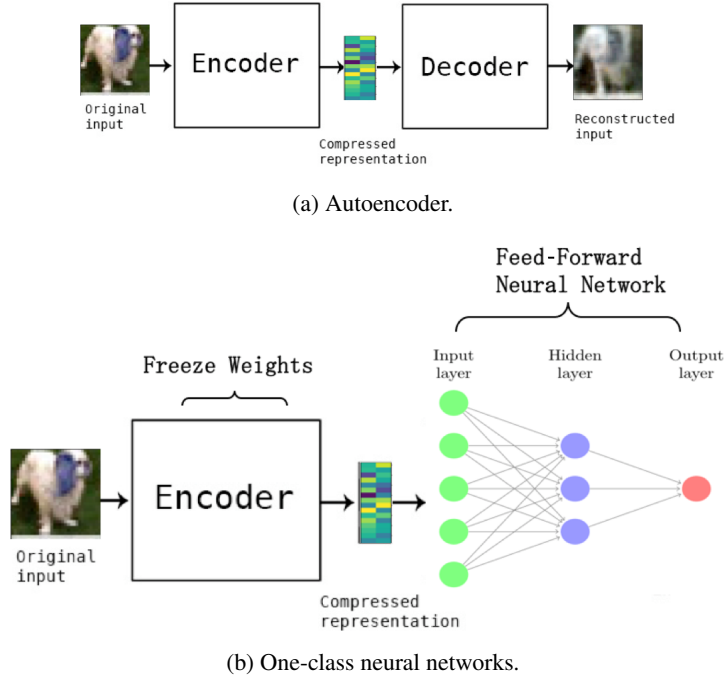
(a) Autoencoder.



(b) One-class neural networks.

Figure 1: Model architecture of Autoencoder and the proposed one-class neural networks (OC-NN).

| Dataset | # Input (features) | # hidden layer( or output) | # optional layer |
|---------|--------------------|-----------------------------|------------------|
| Synthetic | 512 | 128 | 1 |
| MNIST | 32 | 32 | None |
| CIFAR−10 | 128 | 32 | None |
| GTSRB | 128 | 32 | 16 |

Table 2: Summary of best performing feed-forward network architecture's used in OC-NN model for experiments.

### 4.3.2 Deep Baseline Models:

We compare OC−NN models to four deep approaches described Section 4.1. We choose to train DCAE using the Mean sqaure error (MSE) loss since our experiments are on image data. For the DCAE encoder, we employ the same network architectures as we use for Deep SVDD, RCAE, and OC-NN models. The decoder is then constructed symmetrically, where we substitute max-pooling with upsampling. For AnoGAN we follow the implementation as per Radford et al. [2015] and set the latent space dimensionality to 256. For we Deep SVDD, follow the implementation as per Ruff et al. [2018] and employ a two phase learning rate schedule (searching + fine tuning) with initial learning rate $\eta = 10^{-4}$ and subsequently $\eta = 10^{-5}$. For DCAE we train 250 + 100 epochs, for Deep SVDD 150 + 100. Leaky ReLU activations are used with leakiness $\alpha = 0.1$. For RCAE we train the autoencoder using the robust loss and follow the parameter settings as per formulation in Chalapathy et al. [2017].

### 4.3.3 One-class neural Networks (OC-NN):

In OC-NN technique firstly, a deep autoencoder is trained to obtain the representative features of the input as illustrated in Figure 1(a). Then, the encoder layers of this pre-trained autoencoder is copied and fed as input to the feed-forward network with one hidden layer as shown in Figure 1(b). The summary of feed forward network architecture's used for various datasets is presented in Table 2. The weights of encoder network are not frozen (but trained) while we learn feed-forward network weights, following the algorithm summarized in Section 3.3. A feed-forward neural network consisting of single hidden layer, with linear activation functions produced the best results, as per Equation 3. The optimal value of parameter $\nu \in [0, 1]$ which is equivalent to the percentage of anomalies for each data set, is set according to respective outlier proportions.
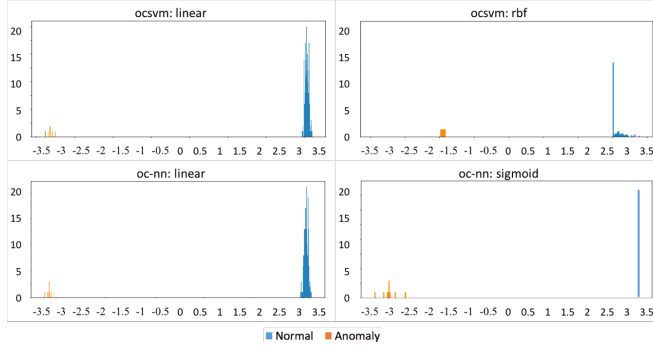
Figure 2: Decision Score Histogram of anomalous vs normal data points, `synthetic` dataset.
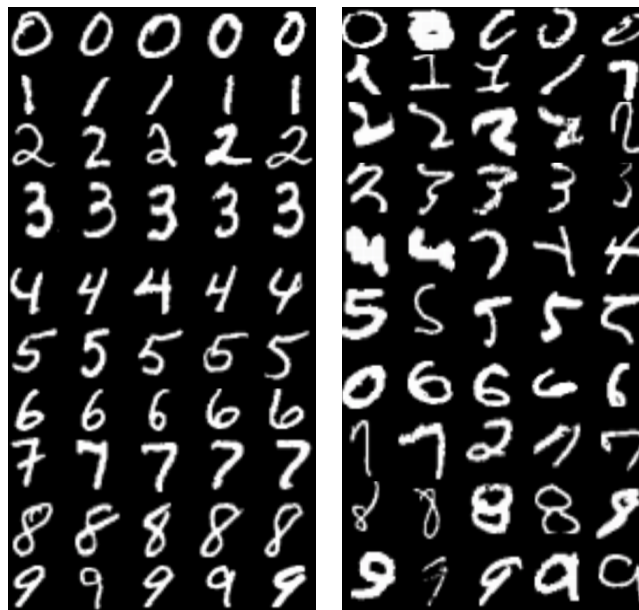
## 5 Experimental Results

In this section, we present empirical results produced by OC-NN model on synthetic and real data sets. We perform a comprehensive set of experiments to demonstrate that on complex data sets OC-NN performs on par with state-of-the-art methods and outperformed conventional shallow methods in some scenarios.

### 5.1 Synthetic Data

Single cluster consisting of 190 data points using $\mu = 0$ and $\sigma = 2$ are generated as normal points , along with 10 anomalous points drawn from normal distribution with $\mu = 0$ and $\sigma = 10$ having dimension $d = 512$. Intuitively, the latter normally distributed points are treated as being anomalous, as the corresponding points have different distribution characteristics to the majority of the training data. Each data point is flattened as a row vector, yielding a $190 \times 512$ training matrix and $10 \times 512$ testing matrix. We use a simple feed-forward neural network model with one-class neural network objective as per Equation 3.
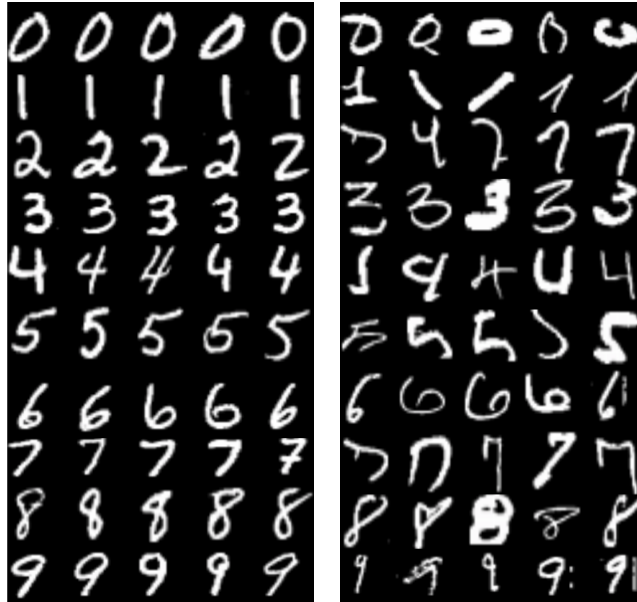
**Results**. From Figure 2, we see that it is a near certainty for all 10 anomalous points, are accurately identified as outliers with decision scores being negative for anomalies. It is evident that, the OC-NN formulation performs on par with classical OC-SVM in detecting the anomalous data points. In the next sections we apply to image and sequential data and demonstrate its effectiveness.



(a) Normal samples.　　　　(b) In-class anomalies.

Figure 3: MNIST Most normal and in-class anomalous MNIST digits detected by RCAE.

| (a) Normal samples. | (b) In-class anomalies. |

Figure 4: MNIST Most normal and in-class anomalous MNIST digits detected by OC-NN.

## 5.2 OC-NN on MNIST and CIFAR-10

**Setup** Both MNIST and CIFAR-10 have ten different classes from which we construct one-class classification dataset. One of the classes is the normal class and samples from the remaining classes represent anomalies. We use the original training and test splits in our experiments and only train with training set examples from the respective normal class. This produces normal instances set of sizes of n $\approx$ 6000 for MNIST and n $\approx$ 5000 for CIFAR-10 respectively. Both MNIST and CIFAR-10 test sets have 10000 samples per class. The training set comprises of normal instances for each class along with anomalies sampled from all the other nine classes. The number of anomalies used within training set consists of 1% normal class for MNIST and 10% of normal class for CIFAR-10 dataset respectively. We pre-process all images with global contrast normalization using the $L1$ norm and finally rescale to [0, 1] via min-max-scaling.

**Network architectures** For both MNIST and CIFAR-10 datasets, we employ LeNet type CNNs, wherein each convolutional module consists of a convolutional layer followed by leaky ReLU activations and $2 \times 2$ max-pooling. On MNIST, we use a CNN with two modules, $8 \times (5 \times 5 \times 1)$-filters followed by $4 \times (5 \times 5 \times 1)$-filters, and a final dense layer of 32 units. On CIFAR-10, we use a CNN with three modules, $32 \times (5 \times 5 \times 3)$-filters, $64 \times (5 \times 5 \times 3)$-filters, and $128 \times (5 \times 5 \times 3)$-filters, followed by a final dense layer of 128 units. We use a batch size of 200 and set the weight decay hyperparameter to $\lambda = 10^{-5}$.

**Results:** Results are presented in Table 3. RCAE clearly outperforms both its shallow and deep competitors on MNIST. On CIFAR-10 the results are convincing but for certain classes the shallow baseline methods outperform the deep models. OC-NN, Soft and One-class Deep SVDD however, shows an overall robust performance. On CIFAR-10 classes such as $AUTOMOBILE$, $BIRD$ and $DEER$ which have less global contrast, as illustrated in Table 3 (indicated in blue color) OC-NN seems to outperform the shallow methods, Soft and One-class Deep SVDD methods, this is indicative of their future potential on similar data instances. It is interesting to note that shallow OCSVM/SVDD and KDE perform better than deep methods on two of the ten CIFAR-10 classes. We can see that normal examples of the classes such as $FROG$ and $TRUCK$ on which OCSVM/SVDD performs best as illustrated in Figure 7 seem to have strong global structures. For example, TRUCK images are mostly divided horizontally into street and sky, and $FROG$ have similar colors globally. For these classes, the performance significantly depends on choice of network architecture. Notably, the One-Class Deep SVDD performs slightly better than its soft-boundary counterpart on both datasets. Figure 3, Figure 5 and Figure 4, Figure 6 show examples of the most normal and most anomalous in-class samples detected by RCAE and OC-NN for MNIST and CIFAR-10 dataset respectively.

(a) Normal samples.

(b) In-class anomalies.

Figure 5: CIFAR-10 Most normal and in-class anomalous CIFAR10 digits detected by RCAE.



(a) Normal samples.

(b) In-class anomalies.

Figure 6: CIFAR-10 Most normal and in-class anomalous CIFAR-10 digits detected by OC-NN.



(a) Normal samples.

(b) In-class anomalies.

Figure 7: Most normal (left) and most anomalous (right) in-class examples determined by OCSVM-SVDD for in which OCSVM-SVDD performs best.

Table 3: Average AUCs in % with StdDevs (over 10 seeds) per method on MNIST and CIFAR-10 dataset.

| NORMAL CLASS | OCSVM / SVDD | KDE | IF | DCAE | ANOGAN | SOFT-BOUND DEEP SVDD | ONE-CLASS DEEP SVDD | OC-NN | RCAE |
|---|---|---|---|---|---|---|---|---|---|
| 0 | $96.75 \pm 0.5$ | $97.1 \pm 0.0$ | $95.32 \pm 1.2$ | $99.90 \pm 0.0$ | $96.6 \pm 1.3$ | $98.66 \pm 1.3$ | $97.78 \pm 0.0$ | $97.60 \pm 1.7$ | $\mathbf{99.92 \pm 0.0}$ |
| 1 | $99.15 \pm 0.4$ | $98.9 \pm 0.0$ | $99.35 \pm 0.0$ | $99.96 \pm 2.1$ | $99.2 \pm 0.6$ | $99.15 \pm 0.0$ | $99.08 \pm 0.0$ | $99.53 \pm 0.0$ | $\mathbf{99.97 \pm 2.2}$ |
| 2 | $79.34 \pm 2.2$ | $79.0 \pm 0.0$ | $73.15 \pm 4.5$ | $96.64 \pm 1.3$ | $85.0 \pm 2.9$ | $88.09 \pm 2.2$ | $88.74 \pm 1.2$ | $87.32 \pm 2.1$ | $\mathbf{98.01 \pm 1.2}$ |
| 3 | $85.88 \pm 1.3$ | $86.2 \pm 0.0$ | $81.34 \pm 2.8$ | $98.42 \pm 0.0$ | $88.7 \pm 2.1$ | $88.93 \pm 3.4$ | $88.26 \pm 3.2$ | $86.52 \pm 3.9$ | $\mathbf{99.25 \pm 0.0}$ |
| 4 | $94.18 \pm 1.5$ | $87.9 \pm 0.0$ | $87.40 \pm 2.6$ | $98.72 \pm 0.0$ | $89.4 \pm 1.3$ | $93.88 \pm 2.3$ | $95.24 \pm 1.4$ | $93.25 \pm 2.4$ | $\mathbf{99.23 \pm 0.0}$ |
| 5 | $72.77 \pm 3.7$ | $73.8 \pm 0.0$ | $73.96 \pm 2.9$ | $97.80 \pm 1.3$ | $88.3 \pm 2.9$ | $84.35 \pm 3.1$ | $83.76 \pm 3.1$ | $86.48 \pm 3.3$ | $\mathbf{99.21 \pm 0.0}$ |
| 6 | $95.14 \pm 1.1$ | $87.6 \pm 0.0$ | $88.65 \pm 0.0$ | $99.74 \pm 0.0$ | $94.7 \pm 2.7$ | $97.74 \pm 0.0$ | $97.99 \pm 0.0$ | $97.12 \pm 1.4$ | $\mathbf{99.81 \pm 1.1}$ |
| 7 | $91.86 \pm 1.6$ | $91.4 \pm 0.0$ | $91.44 \pm 1.8$ | $99.08 \pm 2.2$ | $93.5 \pm 1.8$ | $92.60 \pm 1.2$ | $93.55 \pm 2.3$ | $93.64 \pm 2.1$ | $\mathbf{99.18 \pm 0.0}$ |
| 8 | $88.65 \pm 1.2$ | $79.2 \pm 0.0$ | $75.10 \pm 3.7$ | $96.98 \pm 0.0$ | $84.9 \pm 2.1$ | $90.69 \pm 3.3$ | $90.25 \pm 3.1$ | $88.54 \pm 4.7$ | $\mathbf{98.50 \pm 2.2}$ |
| 9 | $92.53 \pm 1.9$ | $88.2 \pm 0.0$ | $87.59 \pm 1.5$ | $98.04 \pm 1.3$ | $92.4 \pm 1.1$ | $94.28 \pm 2.5$ | $94.12 \pm 2.4$ | $93.54 \pm 3.3$ | $\mathbf{98.98 \pm 1.3}$ |
| AEROPLANE | $60.37 \pm 1.3$ | $61.2 \pm 0.0$ | $63.99 \pm 1.1$ | $71.21 \pm 1.5$ | $67.1 \pm 2.5$ | $66.15 \pm 1.1$ | $67.33 \pm 2.2$ | $60.42 \pm 1.9$ | $\mathbf{72.04 \pm 2.5}$ |
| AUTOMOBILE | $63.03 \pm 1.4$ | $63.0 \pm 0.0$ | $60.56 \pm 1.0$ | $63.05 \pm 2.3$ | $54.7 \pm 3.4$ | $57.64 \pm 3.2$ | $58.14 \pm 3.1$ | $61.97 \pm 2.0$ | $\mathbf{63.08 \pm 2.1}$ |
| BIRD | $63.47 \pm 1.0$ | $50.1 \pm 0.0$ | $64.51 \pm 1.1$ | $71.50 \pm 1.1$ | $52.9 \pm 3.0$ | $61.99 \pm 1.4$ | $61.35 \pm 1.5$ | $63.66 \pm 1.4$ | $\mathbf{71.67 \pm 1.3}$ |
| CAT | $60.25 \pm 0.9$ | $56.4 \pm 0.0$ | $56.16 \pm 2.6$ | $60.57 \pm 0.0$ | $54.5 \pm 1.9$ | $57.56 \pm 4.1$ | $55.72 \pm 1.4$ | $53.57 \pm 2.1$ | $\mathbf{60.63 \pm 1.1}$ |
| DEER | $69.15 \pm 0.8$ | $66.2 \pm 0.0$ | $72.66 \pm 1.1$ | $70.85 \pm 2.2$ | $65.1 \pm 3.2$ | $63.36 \pm 1.3$ | $63.32 \pm 1.2$ | $67.40 \pm 1.7$ | $\mathbf{72.75 \pm 3.3}$ |
| DOG | $66.24 \pm 1.5$ | $62.4 \pm 0.0$ | $61.46 \pm 2.7$ | $62.74 \pm 2.1$ | $60.3 \pm 2.6$ | $58.58 \pm 1.2$ | $58.68 \pm 1.4$ | $56.11 \pm 2.1$ | $\mathbf{63.96 \pm 3.3}$ |
| FROG | $\mathbf{71.57 \pm 1.5}$ | $71.3 \pm 0.0$ | $68.06 \pm 2.1$ | $65.17 \pm 4.1$ | $58.5 \pm 1.4$ | $63.93 \pm 3.1$ | $64.45 \pm 2.1$ | $63.31 \pm 3.0$ | $64.88 \pm 4.2$ |
| HORSE | $63.38 \pm 0.8$ | $62.6 \pm 0.0$ | $63.04 \pm 1.5$ | $61.11 \pm 1.4$ | $62.5 \pm 0.8$ | $60.20 \pm 2.2$ | $59.80 \pm 2.6$ | $60.09 \pm 2.7$ | $\mathbf{63.64 \pm 0.0}$ |
| SHIP | $60.44 \pm 1.1$ | $65.1 \pm 0.0$ | $68.01 \pm 1.3$ | $74.18 \pm 1.2$ | $74.68 \pm 4.1$ | $70.21 \pm 1.1$ | $67.44 \pm 2.2$ | $64.67 \pm 1.6$ | $\mathbf{74.72 \pm 1.1}$ |
| TRUCK | $\mathbf{75.81 \pm 0.8}$ | $74.0 \pm 0.0$ | $72.83 \pm 1.1$ | $71.36 \pm 4.3$ | $66.5 \pm 2.8$ | $72.91 \pm 3.3$ | $68.03 \pm 3.2$ | $60.32 \pm 4.9$ | $74.47 \pm 1.5$ |

Table 4: Average AUCs in % with StdDevs (over 10 seeds) per method on GTSRB stop signs with adversarial attacks.

| METHOD | AUC |
|---|---|
| OC-SVM/SVDD | $52.5 \pm 1.3$ |
| KDE | $51.5 \pm 1.6$ |
| IF | $53.37 \pm 2.9$ |
| AnoGAN | - |
| DCAE | $79.1 \pm 3.0$ |
| SOFT-BOUND. DEEP SVDD | $67.53 \pm 4.7$ |
| ONE-CLASS. DEEP SVDD | $67.08 \pm 4.7$ |
| OC-NN | $63.53 \pm 2.5$ |
| **RCAE** $\lambda = 0$ | $87.45 \pm 3.1$ |
| **RCAE** | $\mathbf{87.39 \pm 2.7}$ |

## 5.3 Detecting Adversarial attacks on GTSRB stop signs using OC-NN

**Setup:** In many applications (eg. autonomous driving) it is of paramount interest to effectively detect the adversarial samples to ensure safety and security. In this experiment, we examine the performance of proposed algorithms on detecting adversarial examples. We consider the "stop sign" class of the German Traffic Sign Recognition Benchmark (GTSRB) dataset, for which we generate adversarial examples from randomly drawn stop sign images of the test set using Boundary Attack Brendel et al. [2017]. We create training dataset consists of n = 1150 examples obtained by combining the normal and anomalous samples in both train and test samples. The number of normal instances n = 1050 stop signs ( 780 from train set + 270 test set ) alongwith 100 adversarial examples are added to obtain training set. We pre-process the data by removing the 10% border around each sign, and then resize every image to $32 \times 32$ pixels following the same setup as in Ruff et al. [2018]. Furthermore, we apply global contrast normalization using the $L1 - norm$ and rescale to the unit interval $[0, 1]$.

**Network architecture** We use a CNN with LeNet architecture having three convolutional modules, $16 \times (5 \times 5 \times 3)$-filters, $32 \times (5 \times 5 \times 3)$-filters, and $64 \times (5 \times 5 \times 3)$-filters, followed by a final dense layer of 32 units. We train with a smaller batch size of 64, due to the dataset size and set again hyperparamter $\lambda = 10^{-6}$.

**Results:** Results Table 4 illustrates the AUC scores obtained. The RCAE outperforms all the other deep models. Figure 8 and Figure 9 shows the most anomalous samples detected by RCAE and OCNN methods respectively, the outliers in this experiment are the images in odd perspectives and the ones that are cropped incorrectly.

## 6 Conclusion

In this paper, we have proposed a one-class neural network (OC-NN) approach for anomaly detection. OC-NN uses a one-class SVM (OC-SVM) like loss function to train a neural network. The advantage of OC-NN is that the features of the hidden layers are constructed for the specific task of anomaly detection. This approach is substantially different from recently proposed hybrid approaches which use deep learning features as input into an anomaly detector. Feature extraction in hybrid approaches is generic and not aware of the anomaly detection task. To learn the parameters of the OC-NN network we have proposed a novel alternating minimization approach and have shown that the optimization of a subproblem in OC-NN is equivalent to a quantile selection

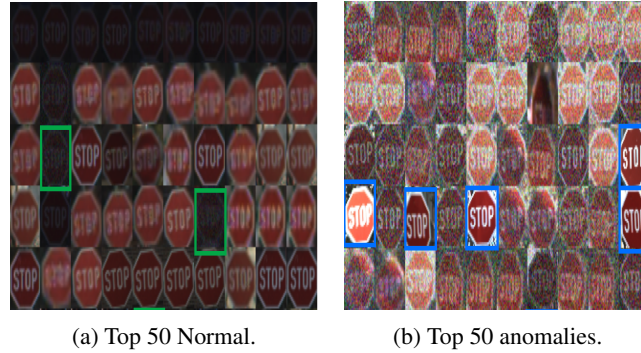(a) Top 50 Normal.      (b) Top 50 anomalies.

Figure 8: Most normal and anomalous stop signs detected by RCAE. Adversarial examples are highlighted in green, Normal samples are highlighted in blue
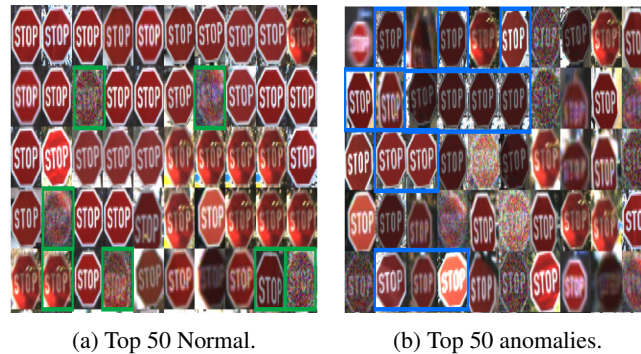


(a) Top 50 Normal.      (b) Top 50 anomalies.

Figure 9: Most normal and anomalous stop signs detected by OC-NN. Adversarial examples are highlighted in green, Normal samples are highlighted in blue

# References

Varun Chandola, Arindam Banerjee, and Vipin Kumar. Outlier detection: A survey. *ACM Computing Surveys*, 2007.

Douglas M Hawkins. *Identification of outliers*, volume 11. Springer, 1980.

Bernhard Schölkopf and Alexander J Smola. Support vector machines, regularization, optimization, and beyond. *MIT Press*, 656:657, 2002.

David MJ Tax and Robert PW Duin. Support vector data description. *Machine learning*, 54(1):45–66, 2004.

Vladimir Naumovich Vapnik and Vlamimir Vapnik. *Statistical learning theory*, volume 1. Wiley New York, 1998.

SVN Vishwanathan, Alexander J Smola, and M Narasimha Murty. Simplesvm. In *Proceedings of the Twentieth International Conference on International Conference on Machine Learning*, pages 760–767. AAAI press, 2003.

Yoshua Bengio, Yann LeCun, et al. Scaling learning algorithms towards ai. *Large-scale kernel machines*, 34(5):1–41, 2007.

Chong Zhou and Randy C Paffenroth. Anomaly detection with robust deep autoencoders. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 665–674. ACM, 2017.

Raghavendra Chalapathy, Aditya Krishna Menon, and Sanjay Chawla. Robust, deep and inductive anomaly detection. *arXiv preprint arXiv:1704.06743*, 2017.

Jerone TA Andrews, Edward J Morton, and Lewis D Griffin. Detecting anomalous data using auto-encoders. *International Journal of Machine Learning and Computing*, 6(1):21, 2016a.

Sarah M Erfani, Sutharshan Rajasegarar, Shanika Karunasekera, and Christopher Leckie. High-dimensional and large-scale anomaly detection using a linear one-class svm with deep learning. *Pattern Recognition*, 58:121–134, 2016.

Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359, 2010.

Jerone TA Andrews, Thomas Tanay, Edward J Morton, and Lewis D Griffin. Transfer representation-learning for anomaly detection. ICML, 2016b.

Charu C Aggarwal. *Outlier Analysis*. Springer, 2nd edition, 2016.

Bernhard Schölkopf, John C Platt, John C Shawe-Taylor, Alex J Smola, and Robert C Williamson. Estimating the support of a high-dimensional distribution. *Neural Computation*, 13(7):1443–1471, Jul 2001.

LJ Cao, Kok Seng Chua, WK Chong, HP Lee, and QM Gu. A comparison of pca, kpca and ica for dimensionality reduction in support vector machine. *Neurocomputing*, 55(1-2):321–336, 2003.

Julia Neumann, Christoph Schnörr, and Gabriele Steidl. Combined svm-based feature selection and classification. *Machine learning*, 61(1-3):129–150, 2005.

Yong Shean Chong and Yong Haur Tay. Abnormal event detection in videos using spatiotemporal autoencoder. In *International Symposium on Neural Networks*, pages 189–196. Springer, 2017.

Erik Marchi, Fabio Vesperini, Stefano Squartini, and Björn Schuller. Deep recurrent neural network-based autoencoders for acoustic novelty detection. *Computational intelligence and neuroscience*, 2017, 2017.

Muhammad Sohaib, Cheol-Hong Kim, and Jong-Myon Kim. A hybrid feature model and deep-learning-based bearing fault diagnosis. *Sensors*, 17(12):2876, 2017.

Lukas Ruff, Robert Vandermeulen, Nico Goernitz, Lucas Deecke, Shoaib Ahmed Siddiqui, Alexander Binder, Emmanuel Müller, and Marius Kloft. Deep one-class classification. In Jennifer Dy and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4393–4402, Stockholmsmässan, Stockholm Sweden, 10–15 Jul 2018. PMLR. URL `http://proceedings.mlr.press/v80/ruff18a.html`.

Huan Xu, Constantine Caramanis, and Sujay Sanghavi. Robust pca via outlier pursuit. In *Advances in Neural Information Processing Systems*, pages 2496–2504, 2010.

Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.

Fei Tony Liu, Kai Ming Ting, and Zhi-Hua Zhou. Isolation forest. In *Data Mining, 2008. ICDM'08. Eighth IEEE International Conference on*, pages 413–422. IEEE, 2008.

Emanuel Parzen. On estimation of a probability density function and mode. *The annals of mathematical statistics*, 33 (3):1065–1076, 1962.

Jonathan Masci, Ueli Meier, Dan Cireşan, and Jürgen Schmidhuber. Stacked convolutional auto-encoders for hierarchical feature extraction. In *International Conference on Artificial Neural Networks*, pages 52–59. Springer, 2011.

Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.

François Chollet et al. Keras. `https://github.com/keras-team/keras`, 2015.

Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.

Yann LeCun, Corinna Cortes, and Christopher JC Burges. Mnist handwritten digit database. *AT&T Labs [Online]. Available: http://yann. lecun. com/exdb/mnist*, 2, 2010.

Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. The german traffic sign recognition benchmark: a multi-class classification competition. In *Neural Networks (IJCNN), The 2011 International Joint Conference on*, pages 1453–1460. IEEE, 2011.

Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, 2009.

Wieland Brendel, Jonas Rauber, and Matthias Bethge. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. *arXiv preprint arXiv:1712.04248*, 2017.