OMASGAN: Out-of-distribution Minimum Anomaly Score GAN for Anomaly Detection

Nikolaos Dionelis, Sotirios A. Tsaftaris, Mehrdad Yaghoobi

SSPD 2022

University of Edinburgh, UK School of Engineering Electronics and Electrical Engineering Institute for Digital Communications (IDCOM)

University Defence Research Collaboration (UDRC) in Signal Processing Track 3 Deep Neural Networks and Machine Learning Methods WP3.1 Robust Generative Neural Networks



https://udrc.eng.ed.ac.uk/track-3-deep-neural-networks-and-machine-learning-methods

1

Contents

- Introduction
 - Anomaly detection/ Out-of-Distribution (OoD) detection in images
- Related Work
 - Generative Adversarial Networks (GAN)
 - Deep generative models knowing what they do not know
- Our proposed OoD Minimum Anomaly Score GAN (OMASGAN) model
 - Main contributions of OMASGAN
 - Active **negative** sampling and training
 - Flowchart diagram; Algorithm
- Evaluation of OMASGAN
 - Results of Leave-One-Out (LOO) evaluation
- Conclusion

Introduction

Anomaly Detection/ OoD Detection in Images

Anomaly detection

- Example: Pedestrians; Van
- Out-of-Distribution (OoD) detection
- Identification of samples that are different from normal data



- <u>Aim</u>: Provide decision mechanisms to decide normal vs abnormal
- Anomalies and abnormal data:
 - Are rare
 - **Might not be known** during training
- <u>Application areas:</u> Important critical systems
 - Security; Safety; Autonomous systems
 - **Defence**; Medical imaging; Healthcare







OMASGAN - Main Contributions

• Problem:

- Models may set **high likelihood** and **low reconstruction loss** to OoD samples
 - Leads to failures to detect anomalies
- Aim: Accurate and robust anomaly detection

• Approach:

- Generative Adversarial Networks (GAN); Generate abnormal samples
- Training with positive and **negative** data; Address: Learning-**OoD-samples** problem

• Main contributions:

- Improved **OoD detection** performance
- Generation of the boundary of the support of the normal class distribution:
 - **OoD minimum-anomaly score** samples
 - Include in training
- Devise a **discriminator** for anomaly detection



Related Work

Related Work and Main Challenges

7

- **Generative models:** Learn the underlying **distribution** of the normal class data
- Problem: Models know what they do not know
 - Set high likelihood and low reconstruction loss to **OoD data** Ο
 - **This decreases** the anomaly detection performance Ο



- Eliminate false negative errors Ο
 - **Failures** to detect anomalies
- Reduce false positives
 - False alarms of anomalies



Related Research on Active Negative Training

- Old is Gold (OGNet) [1]: Uses a pseudo-anomaly module to produce OoD samples
 - Restrictive definition of anomaly
 - Blurry reconstructions before convergence
 - Single-epoch reconstructed images
- Active negative training but not active negative sampling
 - Does not cover the **OoD part** of the data space
- Discriminator for AD: Separates good from bad quality reconstructions



[1] M. Zaheer, et al. Old is Gold: Redefining the Adversarially Learned One-Class Classifier Training Paradigm. CVPR 2020

Related Work on Anomaly Detection

- Anomaly detection using the reconstruction error
 - Unsupervised training for OoD detection: No class labels
 - GAN for Anomaly detection (GANomaly) [2]: Encoder-Decoder-Encoder
- Fence GAN [3], Minimum Likelihood GAN [4]:
 - Inference: Discriminator
- Boundary of Data Distribution Support Generator (BDSG) [5]:
 - Flow-based invertible generative models
 - Invertible Residual Networks



[2] S. Akcay, et al., "GANomaly: Semi-Supervised Anomaly Detection via Adversarial Training," in Proc. ACCV, 2018
[3] P. Cuong Ngo, A. A. Winarto, C. K. Li Kou, S. Park, F. Akram, H. K. Lee. "Fence GAN: Towards Better Anomaly Detection", In Proc. 31st International Conference on Tools with Artificial Intelligence (ICTAI), 2019
[4] C. Wang, Y. Zhang, and C. Liu, "Anomaly Detection via Minimum Likelihood Generative Adversarial Networks". In Proc. 24th International Conference on Pattern Recognition (ICPR), 2018

[5] N. Dionelis, M. Yaghoobi, and S. A. Tsaftaris, "Boundary of Distribution Support Generator (**BDSG**): Sample Generation on the Boundary," in Proc. IEEE International Conference on Image Processing (ICIP), 2020

OoD Minimum Anomaly Score GAN (OMASGAN)

OoD Minimum Anomaly Score GAN Model

- Main problem:
 - Models may set **high likelihood** and **low reconstruction loss** to OoD samples
 - Overall decreases Anomaly Detection (AD) performance
- AD models underperform due to the rarity of anomalies
- We develop: The proposed OoD Minimum Anomaly Score GAN (OMASGAN) model
 - Self-generate abnormal/ OoD samples
 - On the boundary of the support of the normal class data distribution
 - Invertibility or probability density: Not needed
 - Impose low likelihood on this learned data distribution boundary
 - Include the proposed OoD minimumanomaly score samples in training
 - Devise a discriminator for AD trained on both negative and positive data



Contributions: OMASGAN Methodology

- Focus on: GANs to improve anomaly detection performance
- GANs: Able to learn complex data
 - Achieve convergence in distribution metrics
- Our proposed methodology: Use only samples from the normal class
 - Generate the **OoD boundary** of the support of the **data distribution**
- Find the **OoD** minimum-anomaly-score samples on the distribution boundary
 - With I_p-norm and dispersion regularisation
 - Perform self-supervised learning
- Integrate with

any **f-divergence GAN** generative model



"OMASGAN: Out-of-distribution Minimum Anomaly Score GAN for Anomaly Detection," Nikolaos Dionelis, Sotirios A. Tsaftaris, and Mehrdad Yaghoobi, SSPD 2022

OMASGAN Flowchart



- Train a **f-divergence GAN**: Obtain the **generator**, G(**z**)
- Train a **boundary data generator**, B(**z**), to obtain **OoD** distribution boundary samples
 - Active negative sampling
- OMASGAN generates OoD data, B(z)
 - Incorporate OoD Minimum Anomaly Score (OMAS) samples in training
- Active negative training of our proposed generator, G'(z)
- Devise a **discriminator** for anomaly detection

Our Proposed OMASGAN Algorithm

OMASGAN - Our Proposed Algorithm

- Train a f-divergence GAN to establish a distribution metric
 - $\circ \min_{G} \max_{D} E_{\mathbf{x}} \log(D(\mathbf{x})) + E_{\mathbf{z}} \log(1 D(G(\mathbf{z})))$
 - where $z \sim p_z(z)$, $G(z) \sim p_g(x)$, $x \sim p_x(x)$
- Train the proposed distribution boundary model, B(z)
 - Perform active negative sampling
 - $\circ \min_{\mathsf{B}} -\mathsf{m}(\mathsf{B}(\mathsf{z}), \, \mathsf{G}(\mathsf{z})) + \lambda \, \mathsf{d}(\mathsf{B}(\mathsf{z}), \, \mathsf{G}(\mathsf{z})) + \mu \, \mathsf{s}(\mathsf{B}(\mathsf{z}), \, \mathsf{z})$
 - where m is a f-divergence metric
- Active negative training: More accurately learn the underlying data distribution
 - Train G'(z) to address the learning-OoD-samples problem of G(z)
 - $\circ \min_{G'} \max_{C} \alpha E_{\mathbf{x}} \log(C(\mathbf{x})) + \beta E_{\mathbf{z}} \log(1-C(G'(\mathbf{z}))) + (1-\beta) E_{\mathbf{z}} \log(1-C(B(\mathbf{z}))) + (1-\alpha) E_{\mathbf{z}} \log(C(G(\mathbf{z})))$
 - where C is a discriminator
- Train a discriminator, J, to perform active negative learning for AD
 - $\circ \max_{J} \gamma E_{\mathbf{x}} \log(1 J(\mathbf{x})) + (1 \gamma) E_{\mathbf{z}} \log(1 J(G'(\mathbf{z}))) + E_{\mathbf{z}} \log(J(B(\mathbf{z})))$
 - J learns to **separate** B(**z**) from G'(**z**) and **x**

Generation of OMAS Samples

- Generate **OoD minimum-anomaly-score** samples, B(z)
 - On the **boundary of the normal class data distribution**
 - $\circ \quad \text{arg min}_{\text{B}} \text{-m}(\text{B}(\textbf{z}), \text{ G}(\textbf{z})) + \lambda \text{ d}(\text{B}(\textbf{z}), \text{ G}(\textbf{z})) + \mu \text{ s}(\text{B}(\textbf{z}), \textbf{z})$

$$d(B(\mathbf{z};\boldsymbol{\theta_b}), G(\mathbf{z})) = \min_{j=1,\dots,Q} ||B(\mathbf{z};\boldsymbol{\theta_b}) - G(\mathbf{z}_j)||_p^q$$
$$s(B(\mathbf{z}_i;\boldsymbol{\theta_b}), \mathbf{z}_i) = \frac{1}{N-1} \sum_{j=1, j \neq i}^N \frac{||\mathbf{z}_i - \mathbf{z}_j||_p^q}{||B(\mathbf{z}_i;\boldsymbol{\theta_b}) - B(\mathbf{z}_j;\boldsymbol{\theta_b})||_p^q}$$

- Train our proposed GAN-based boundary formation model, B(z), for AD:
 - Propose a loss with three terms
- B(z; θ); Run Gradient Descent on the proposed loss; Obtain θ
 - Batch size, N
 - Inference sample size, Q, for G(z)

- Force the samples to the boundary of the data distribution
- Effectively address mode collapse
- First term: Decreasing function of a distribution divergence metric
- Second term: I_p-norm distance
- Third term: Scattering, dispersion
 - Capture all the modes

OMASGAN Inference Mechanism

- Anomaly score/ OoD score based on the Anomaly Discriminator, J, and the f-divergence distribution metric
 - The **discriminator**, J, is trained to:
 - Separate the normal class distribution from its complement
- **f-divergence distribution metric:** Used for training and during **inference**
- **f-divergence** for probability distributions **P** and **Q**: **fD(P, Q)**
- For a queried test sample, x*:
 - **Calculate** fD(G', δ_{x^*})
 - where δ_{x^*} is a Dirac function centered at x^*
- Our proposed anomaly/ OoD score: $AS(\mathbf{x}^*; J, G') = J(\mathbf{x}^*) + \lambda fD(G', \delta_{\mathbf{x}^*})$
- Classification decision: **x*** is from the **normal class** if AS(**x***; J, G') < τ
 - \circ where τ is a threshold
 - **x*** is **abnormal** otherwise

Evaluation of OMASGAN

AD Evaluation of the Proposed Model

- Evaluation methods for anomaly/ OoD detection:
 - Leave-one-out (LOO) evaluation
 - Normal class: 9 classes from a benchmark dataset with 10 classes
 - Abnormal class: Left-out class
- Evaluation metrics: Algorithm convergence criteria; Area Under the Receiver Operating Characteristics Curve (AUROC)
- Examined datasets:
 - Synthetic; MNIST; CIFAR-10
- Baselines: GANomaly [2]; VAE EGBAD; AnoGAN; FenceGAN [3]; MinLGAN [4] BDSG [5]; TailGAN





Our proposed model:

- OMASGAN
 - KLWGAN-based OMASGAN
 - f-GAN-based OMASGAN

Evaluation of the OMASGAN Model on MNIST

- Performance of OMASGAN on MNIST data in AUROC
 - Compared to GAN and AE baselines using LOO evaluation Ο



OMASGAN = BDSG = TailGAN = EGBAD = AnoGAN = GANomaly = VAE

Evaluation of OMASGAN on CIFAR-10 Data

- Performance of OMASGAN in AUROC on the CIFAR-10 dataset
 - Comparison with GAN and AE baselines using LOO evaluation



Evaluation of OMASGAN - Ablation Study



- OMASGAN • Task3 • Task1

- Ablation study of OMASGAN in AUROC on MNIST (left) and on CIFAR-10 (bottom)
- We examine the impact of our loss functions
 - Using LOO evaluation

- **OMASGAN** outperforms:
 - Our chosen base model, KLWGAN
 - Ablation study: G'(z)



Conclusion

Conclusion

- OMASGAN: GANs to improve the anomaly detection/ OoD detection performance
- Use only data from the normal class
 - Perform active negative sampling
 - Generate abnormal distribution boundary samples
 - Perform active negative training for anomaly detection
- Address the learning-OoD-samples problem of generators
 - Contrastive negative training to alleviate the problem of deep generative models knowing what they do not know (setting high likelihood to OoD data)
- Effectively tackle the rarity of anomalies problem
- Propose a discriminator-based anomaly score/ OoD score
- **OMASGAN outperforms** the examined baselines
 - Using the LOO evaluation methodology
 - In **AUROC**, on the MNIST and CIFAR-10 datasets

Thank you very much for your attention!

Contact email: Nikolaos.Dionelis@ed.ac.uk

References

- [1] N. Dionelis, M. Yaghoobi, and S. A. Tsaftaris. "Boundary of Distribution Support Generator (**BDSG**): Sample Generation on the Boundary", in Proc. IEEE International Conference on Image Processing (ICIP), 2020
- [2] N. Dionelis, S. A. Tsaftaris, and M. Yaghoobi. OMASGAN: Out-of-Distribution Minimum Anomaly Score GAN for Sample Generation on the Boundary, <u>https://github.com/nd1511/OMASGAN_</u>
- [3] N. Dionelis, M. Yaghoobi, and S. A. Tsaftaris. "Tail of Distribution GAN (**TailGAN**): Generative Adversarial Network Based Boundary Formation", in Proc. Sensor Signal Processing for Defence (SSPD), 2020
- [4] S. Akcay, A. Atapour-Abarghouei, and T. Breckon. GANomaly: Semi-Supervised Anomaly Detection via Adversarial Training. arXiv preprint, arXiv:1805.06725 [cs.CV], 2018
- [5] H. Zenati, C. Foo, B. Lecouat, G. Manek, and V. Ramaseshan Chandrasekhar. Efficient GAN-Based Anomaly Detection. Workshop Track in International Conference on Learning Representations (ICLR), 2018
- [6] T. Schlegl, P. Seeböck, S. Waldstein, U. Schmidt-Erfurth, and G. Langs. Unsupervised Anomaly Detection with GANs to Guide Marker Discovery. In Proc. Information Processing in Medical Imaging (IPMI), 2017
- [7] P. Cuong Ngo, A. A. Winarto, C. K. Li Kou, S. Park, F. Akram, H. K. Lee. "Fence GAN: Towards Better Anomaly Detection", In Proc. 31st International Conference on Tools with Artificial Intelligence (ICTAI), 2019
- [8] C. Wang, Y. Zhang, and C. Liu, "Anomaly Detection via Minimum Likelihood Generative Adversarial Networks". In Proc. 24th International Conference on Pattern Recognition (ICPR), 2018
- [9] R. Devon Hjelm, A. Paul Jacob, T. Che, A. Trischler, K. Cho, and Y. Bengio. Boundary-Seeking Generative Adversarial Networks. arXiv preprint, arXiv:1702.08431 [stat.ML], 2018
- [10] L. Ruff, R. Vandermeulen, N. Gornitz, L. Deecke, S. Siddiqui, A. Binder, E. Muller, and M. Kloft. Deep One-Class Classification. In Proc. International Conference on Machine Learning (ICML), 2018
- [11] L. Ruff, R. A. Vandermeulen, et al., "Deep Semi-Supervised Anomaly Detection", arXiv:1906.02694, 2020