

Research Statement

PANG Guansong

School of Computing and Information Systems, Singapore Management University

Tel: (65) 6828-4864; Email: gspang@smu.edu.sg

22 December 2023

Background

One focused theme of my research is to handle rare, abnormal, or unknown instances in large-scale data. This plays a vital role in a broad range of domains, such as preventing the loss of billions of dollars by their application to fraud detection and anti-money laundering in fintech, saving lives through early disease detection, safeguarding large-scale computer networks and data centers from malicious attacks by their use in intrusion detection, equipping AI systems with the ability to work safely in open worlds, and accelerating scientific discovery by their application to identify novel observations in large-scale medicine, physics, chemistry and material science data. To address this general problem, I have been dedicating to pushing the boundary of various crucial related tasks, such as anomaly detection and robust deep learning (robustness *w.r.t.* out-of-distribution samples, long-tail distribution, and adversarial examples).

I have been also very keen on representation learning that aims at learning expressive feature representations of data samples in different form, such as image, video, graph, and sequence data. The success of most machine learning algorithms depends on feature representation, which is especially true for real-world problems and applications where the data contain intricate dependency and structure. To effectively learn from those data, one way is feature engineering that focuses on manually constructing a set of features to enable the subsequent learning algorithms. Representation learning instead extracts expressive features from the data without human engineering efforts, which can result in more optimal feature representations while at the same time being significantly cost-effective in human involvement.

Research Areas

Anomaly Detection

Deep learning for anomaly detection.

Deep learning can substantially enhance traditional anomaly detection in many ways, e.g., learning discriminative, non-linear features in an end-to-end fashion, anomaly score optimization, etc. It is an exciting emerging research direction in anomaly detection due to the remarkable performance of deep anomaly detection on different types of datasets. We have explored several important directions to devise deep anomaly detectors. One way is to devise objective functions to learn feature representations that are optimized for some existing state-of-the-art shallow anomaly detection measures, which allows us to make use of the plethora of the existing shallow methods. We introduce such an objective function to enable traditional distance-based anomaly measure in [KDD'18], in which we show that the tailored feature representations learned by our method can significantly improve the performance of the shallow anomaly measures than working on the original feature space, while at the same time largely speeding up the online detection stage. We have also explored another important approach that focuses on devising new anomaly detection approaches based on deep learning techniques, including new loss functions [KDD'19; CVPR'20; AAI'22; KDD'23] and new detection models [IJCAI'20; ICCV'21; WSDM'22; TKDE'23]. In [MICCAI'21; MedIA'23], we also introduce new self-supervised pre-training methods to enable deep anomaly detection on datasets with small training samples. Despite remarkable progress brought by deep anomaly detection, there are still a number of largely unsolved challenges in the area,

such as effective supervisory signals for training deep anomaly detectors, the detection of hard anomalies, robustness to distribution shift [ICCV'23], etc. We will continue dedicating to this area to address those specific challenges.

Non-i.i.d. anomaly detection.

'Data is independent and identically distributed (*i.i.d.*)' is a widely used fundamental assumption made in most machine learning (including anomaly detection) algorithms and theories. However, this assumption is violated in many real-world applications. We explore and model the coupling relations between the outlier factors of different entities, such as feature values, features, and data instances to learn anomaly scores that well capture the underlying intricate abnormal behaviors. A large number of couplings relations are exploited, such as conditional cascade [IJCAI'16], binary cascade [IJCAI'17], or high-order cascade of outlier factors [DMKD'21] in modeling the anomalousness due to feature value interactions, pairwise interactions between feature-level outlier factors [ICDM'16], and sequential couplings of the anomaly scores of data instances [AAAI'18]. The resulting models enable significantly improved performance in the detection of subtle anomalies, or anomalies in high-dimensional data, in which anomalies are masked when not considering these non-i.i.d. data characteristics. In addition to the non-independent aspect, it is also important to consider the non-identically-distributed aspect, as we showed in [TKDD'20], the distribution of the abnormal behaviors presented in one feature differs largely from that in the other features. Simple models that consider this feature heterogeneity can perform much better than complex models that ignore such heterogeneities. Studies in this direction often require some prior knowledge of or assumptions on the underlying coupling relations in the datasets, leading to models that may be not generalizable to other datasets. It is important to develop tools for more generalizable non-i.i.d. detectors.

Graph data is typical non-i.i.d. data. In recent years, we utilize the power of graph neural networks and devise various effective optimization objectives to support the training of GNN-based anomaly detectors, including anomalous node detection [SDM'23; AAAI'23; NeurIPS'23] and anomalous graph detection [WSDM'22; ECMLPKDD'23; DSAA'23]. Graph data is ubiquitous in many application domains, where there are various largely unsolved challenges, such as the lack of tailored optimization objectives for graph anomaly detection and large labeled training samples, the heterogeneity in the graph, etc. We are committed to addressing these challenges in the coming few years to make our models ready for deployment in real-world applications.

Open-world Machine Learning

Out-of-distribution detection.

Detecting out-of-distribution (OOD) samples is crucial to the safe deployment of deep learning systems in real-world applications. OOD detection can be considered as a special case of anomaly detection, where we are given normal class data to train robust deep learning systems that detect novel samples drawn from outside the training distribution, while at the same time accurately classifying samples into the normal classes. This problem exists in a wide range of tasks, such as image classification, semantic segmentation, and object detection. In [ECCV'22], we introduce a novel approach, named PEBAL, that learns an extra pixel-level class (i.e., outlier class) for the unknowns via abstention learning for detecting OOD pixels in semantic segmentation. The approach enforces a small penalty to the ML model if it abstains from classifying unknown/uncertain instances into the anomaly class; whereas a large penalty is enforced if it abstains from classifying known instances. The approach shows promising results in autonomous driving urban scenes and is among the best performers on different datasets of the public leaderboard for autonomous driving semantic segmentation¹. In [ICCV'23b], we introduce an improved method over PEBAL that

¹ <https://segmentmeifyoucan.com/leaderboard>

can effectively maintain the in-distribution segmentation accuracy, while largely enhancing the OOD detection performance. It is also robust to distribution shift from city-view scenes to country-view driving scenes. We are working on tackling this problem in a variety of other realistic application settings.

Open-set recognition.

Existing anomaly detection methods overwhelmingly focus on training detection models using exclusively normal data (semi-supervised anomaly detection paradigm) or unlabeled data (unsupervised anomaly detection paradigm). These models are not fed with any labelled anomaly data, and thus, they lack knowledge about genuine anomalies, learning features that are not discriminative to distinguish some anomalies from normal data. We promote a supervised anomaly detection paradigm, open-set supervised anomaly detection, that utilizes a small number of labeled anomaly examples to learn anomaly-informed detection models to address this issue. These few-shot anomaly examples may originally come from a deployed detection system, e.g., a few successfully detected defects, credit card frauds, or network intrusion samples; they may be from users/human experts, such as a small number of bank frauds or lesion images that are reported or confirmed by users/human experts. This paradigm offers significantly improved detection accuracy by utilizing few-shot labelled anomaly samples (and large unlabeled) with trivial/small human annotation cost. In [KDD'18], we introduce an anomaly query neural network that learns a new feature representation space so that the few anomaly examples have larger distance-based anomaly scores than that of pseudo normal instances in the new feature space. We further introduce a new loss function, named deviation loss [KDD'19], to achieve more sample-efficient and optimal learning of anomaly scores. These two methods focus on the exploitation of the labeled anomaly data, ignoring valuable information hidden in large-scale unlabeled data where most data are normal instances. In [KDD'21], we introduce a deep reinforcement learning approach to leverage the limited labeled anomaly examples, while at the same time actively exploring the supervisory signals from the large unlabeled data. These supervised detection models often show significantly improved performance on detecting anomalies that are similar to the seen anomalies during training, but they may become less effective in detecting unseen anomalies than the unsupervised detectors. Generalizing to both seen and unseen anomalies is generally required in real-world applications. We are working on different methods to achieve this goal. One of our recent achievements in this research line is [CVPR'22], where we introduce the task of open-set supervised anomaly detection and propose a novel approach that learns disentangled representations of three general types of abnormalities, illustrated by few-shot instances of known anomaly classes, pseudo anomaly classes, and those that largely deviates from normal instances. By doing so, we can effectively detect both known and unknown anomaly classes. In [CVPR'23], we explore a relatively new area, open-set few-shot recognition, and introduce a new state-of-the-art method for open-set recognition with only few-shot training samples.

Representation Learning

Fine-grained representation learning.

Representation learning is one of the key driving forces to the tremendous success of deep learning across different application domains. There have been numerous studies dedicated to address different challenges in this area. One major challenge is how to learn expressive representations in datasets where the intra-class holistic features are large while the inter-class holistic features are small; the original features are discriminative only when looking at some local fine-grained features. This issue presents significant challenges to popular feature learning techniques, such as some popular loss functions (e.g., triplet loss, contrastive loss, etc.), that focus on learning holistic features. To address this issue, we introduce a novel fine-grained difference-aware loss function in [TMM'21]. The proposed loss function can substantially enhance the capability of learning fine-grained discriminative features, especially in distinguishing some hard examples that are distant from their genuine classes

while being close to other classes due to the background features. The effectiveness of this loss function was justified in the person re-identification (ReID) task. This issue is particularly crucial in occluded ReID where the targeted persons are occluded by some unknown objects [ICCV'21b], or in bird-view ReID where only very limited appearance features of persons are visible due to the large vertical angle between the camera and the persons [ICCV'21a]. In [ICCV'21b], we explore the combination of occlusion-based data augmentation and a simpler fine-grained difference-aware loss function to learn the discriminative features from the scattered, small non-occluded body parts. In [ICCV'21a], we introduce a multi-scale cross-attention framework to tackle the issue. We plan to explore more advanced loss functions and network architectures to learn more discriminative features to further largely reduce the high false positives/negatives in existing models.

Self-supervised representation learning.

In many real-world applications, it is very costly to obtain large-scale labeled data, and exclusively fitting the labeled data is prone to overfitting. To address these issues, self-supervised feature learning has been emerging as a very popular research line, in which we explore supervisory signals embedded in the data itself, or by some simple data augmentation techniques, such as to predict the relative position of an image patch w.r.t. other image patches, to predict whether a given sample is an augmented sample of itself or it is a different sample, etc., to train learning models. By doing so, this self-supervised learning of features produces feature representations that contain rich semantics supervised learning tasks may not be able to learn, helping complement and regularize the supervised learning. In [IJCAI'20], we introduce a novel self-supervised learning task, to predict the pairwise distance based on features derived from a randomly initialized neural network. We show the method can learn good manifold and similarity information, largely lifting the performance of existing clustering and anomaly detection techniques. In [WSDM'22], we introduce a self-supervised graph representation learning method, which learns expressive node-level and graph-level representations by enforcing a graph neural network (GNN) model to distill knowledge from a randomly initialized GNN at both the node and graph levels. In [INS'22], a self-supervised learning approach based on neighborhood similarity is introduced to enable the learning of Euclidean and hyperbolic graph representations. There are many existing opportunities in this research area. We plan to continue our efforts on challenging datasets, such as graph data, video/image data, and time series data.

Machine Learning for Cybersecurity

Machine learning has been intensively explored in the cybersecurity domain. We have been exploring the use of machine learning techniques such as anomaly detection to empower cybersecurity applications. We have been testing our anomaly detection models on datasets in different cybersecurity related contexts, including intrusion detection [KDD'18; KDD'19; KDD'21], malicious URL detection [KDD'18; KDD'19], web spam detection [KDD'18]. Particularly, we are very interested in the detection of zero-day network attacks, i.e., attacks that are unknown and not revealed by any researchers or organization. The key challenge here is how we can learn from the known network attacks and generalize the learned attack patterns to detect these unknown attacks. We show in our recent work [KDD'21] that we can use reinforcement learning models to effectively learn patterns of suspicious unknown attacks from large-scale unlabeled data. However, this work assumes that the unlabeled data contains the zero-day attack data, which may therefore fail when this assumption does not hold. We are exploring more advanced tools that enable the detection of this type of attacks without requiring their presence in any form in our training data.

Smart Healthcare

Many problems in healthcare have been formulated as a binary (e.g., benign vs malignant) classification task. Unfortunately, this formulation works ineffectively when the abnormal class

exhibits irregularly variant features or distribution shift from known abnormalities. This can lead to the misclassification of malignant cases and ultimately catastrophic loss. Our research finds that anomaly detection-based approaches that focus on modeling the normal class only can perform significantly better than the classification-based approaches in such cases. This is verified by our extensive experiments in prediction of COVID-19 cases using chest X-ray images [TMI'20], malignant lesion detection in several other organs [MICCAI'21; MICCAI'22; MedIA'23], and prediction of children depression using multivariate temporal data [PAKDD'22]. It would be important to explore this intuition in healthcare application areas where we may have a large number of normal (sub)classes and abnormal classes.

Selected Publications and Outputs

- [IJCAI'16] Pang, G., Cao, L., & Chen, L. (2016, January). Outlier detection in complex categorical data by modelling the feature value couplings. In *IJCAI International Joint Conference on Artificial Intelligence*.
- [ICDM'16] Pang, G., Cao, L., Chen, L., & Liu, H. (2016, December). Unsupervised feature selection for outlier detection by modelling hierarchical value-feature couplings. In *2016 IEEE 16th International Conference on Data Mining (ICDM)* (pp. 410-419). IEEE.
- [IJCAI'17] Pang, G., Cao, L., Chen, L., & Liu, H. (2017, January). Learning homophily couplings from non-IID data for joint feature selection and noise-resilient outlier detection. In *IJCAI International Joint Conference on Artificial Intelligence*.
- [AAAI'18] Pang, G., Cao, L., Chen, L., Lian, D., & Liu, H. (2018, April). Sparse modeling-based sequential ensemble learning for effective outlier detection in high-dimensional numeric data. In *Thirty-second AAAI conference on artificial intelligence*.
- [KDD'18] Pang, G., Cao, L., Chen, L., & Liu, H. (2018, July). Learning representations of ultrahigh-dimensional data for random distance-based outlier detection. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining* (pp. 2041-2050).
- [ACMMM'19] Yan, C., Pang, G., Bai, X., Shen, C., Zhou, J., & Hancock, E. (2019, October). Deep hashing by discriminating hard examples. In *Proceedings of the 27th ACM International Conference on Multimedia* (pp. 1535-1542).
- [KDD'19] Pang, G., Shen, C., & van den Hengel, A. (2019, July). Deep anomaly detection with deviation networks. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining* (pp. 353-362).
- [IJCAI'20] Wang, H., Pang, G., Shen, C., & Ma, C. (2019). Unsupervised representation learning by predicting random distances. In *IJCAI International Joint Conference on Artificial Intelligence*.
- [CVPR'20] Pang, G., Yan, C., Shen, C., Hengel, A. V. D., & Bai, X. (2020). Self-trained deep ordinal regression for end-to-end video anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 12173-12182).
- [TMI'20] Zhang, J., Xie, Y., Pang, G., Liao, Z., Verjans, J., Li, W., ... & Xia, Y. (2020). Viral pneumonia screening on chest x-rays using Confidence-Aware anomaly detection. *IEEE transactions on medical imaging*, 40(3), 879-890.
- [TKDD'20] Pang, G., & Cao, L. (2020). Heterogeneous univariate outlier ensembles in multidimensional data. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 14(6), 1-27.
- [DMKD'21] Pang, G., Cao, L., & Chen, L. (2021). Homophily outlier detection in non-IID categorical data. *Data Mining and Knowledge Discovery*, 1-62.
- [TMM'21] Yan, C., Pang, G., Bai, X., Liu, C., Xin, N., Gu, L., & Zhou, J. (2021). Beyond triplet loss: person re-identification with fine-grained difference-aware pairwise loss. *IEEE Transactions on Multimedia*.
- [ICCV'21a] Yan, C., Pang, G., Wang, L., Jiao, J., Feng, X., Shen, C., & Li, J. (2021). BV-Person: A Large-Scale Dataset for Bird-View Person Re-Identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 10943-10952).
- [ICCV'21b] Yan, C., Pang, G., Jiao, J., Bai, X., Feng, X., & Shen, C. (2021). Occluded Person Re-Identification With Single-Scale Global Representations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 11875-11884).
- [ICCV'21c] Tian, Y., Pang, G., Chen, Y., Singh, R., Verjans, J. W., & Carneiro, G. (2021). Weakly-supervised Video Anomaly Detection with Robust Temporal Feature Magnitude Learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*.

- [KDD'21] Pang, G., van den Hengel, A., Shen, C., & Cao, L. (2021, August). Toward deep supervised anomaly detection: Reinforcement learning from partially labeled anomaly data. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining* (pp. 1298-1308).
- [MICCAI'21] Tian, Y., Pang, G., Liu, F., Chen, Y., Shin, S. H., Verjans, J. W., ... & Carneiro, G. (2021, September). Constrained contrastive distribution learning for unsupervised anomaly detection and localisation in medical images. In *International Conference on Medical Image Computing and Computer-Assisted Intervention* (pp. 128-140). Springer, Cham.
- [WSDM'22] Rongrong Ma, Guansong Pang, Ling Chen, and Anton van den Hengel. (2022). Deep Graph-level Anomaly Detection by Glocal Knowledge Distillation. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*.
- [ECCV'22] Tian, Y., Liu, Y., Pang, G., Liu, F., Chen, Y., & Carneiro, G. (2022). Pixel-wise energy-biased abstention learning for anomaly segmentation on complex urban driving scenes. In *European Conference on Computer Vision* (pp. 246-263). Springer, Cham.
- [CVPR'22] Ding, C., Pang, G., & Shen, C. (2022). Catching Both Gray and Black Swans: Open-set Supervised Anomaly Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 7388-7398).
- [AAAI'22] Chen, Y., Tian, Y., Pang, G., & Carneiro, G. (2022, June). Deep one-class classification via interpolated gaussian descriptor. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 36, No. 1, pp. 383-392).
- [INS'22] Xu, X., Pang, G., Wu, D., & Shang, M. (2022). Joint hyperbolic and Euclidean geometry contrastive graph neural networks. *Information Sciences*, 609, 799-815.
- [MICCAI'22] Tian, Y., Pang, G., Liu, F., Liu, Y., Wang, C., Chen, Y., ... & Carneiro, G. (2022). Contrastive transformer-based multiple instance learning for weakly supervised polyp frame detection. In *International Conference on Medical Image Computing and Computer-Assisted Intervention* (pp. 88-98). Springer, Cham.
- [PAKDD'22] Pang, G., Pham, N. T. A., Baker, E., Bentley, R., & van den Hengel, A. (2022). Deep Depression Prediction on Longitudinal Data via Joint Anomaly Ranking and Classification. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining* (pp. 236-248). Springer, Cham.
- [ICCV'23] Tri Cao, Jiawen Zhu, and Guansong Pang. "Anomaly Detection under Distribution Shift." In: ICCV'23.
- [KDD'23] Pang, Guansong, Chunhua Shen, Huidong Jin, and Anton van den Hengel. "Deep weakly-supervised anomaly detection." In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 1795-1807. 2023.
- [MedIA23] Tian, Yu, Fengbei Liu, Guansong Pang, Yuanhong Chen, Yuyuan Liu, Johan W. Verjans, Rajvinder Singh, and Gustavo Carneiro. "Self-supervised pseudo multi-class pre-training for unsupervised anomaly detection and segmentation in medical images." *Medical image analysis* 90 (2023): 102930.
- [TKDE'23] Xu, Hongzuo, Guansong Pang, Yijie Wang, and Yongjun Wang. "Deep isolation forest for anomaly detection." *IEEE Transactions on Knowledge and Data Engineering* (2023).
- [ICCV'23b] Liu, Yuyuan, Choubo Ding, Yu Tian, Guansong Pang, Vasileios Belagiannis, Ian Reid, and Gustavo Carneiro. "Residual pattern learning for pixel-wise out-of-distribution detection in semantic segmentation." In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 1151-1161. 2023.
- [AAAI'23] Wang, Qizhou, Guansong Pang, Mahsa Salehi, Wray Buntine, and Christopher Leckie. "Cross-domain graph anomaly detection via anomaly-aware contrastive alignment." In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 4, pp. 4676-4684. 2023.
- [ECMLPKDD'23] Niu, Chaoxi, Guansong Pang, and Ling Chen. "Graph-level anomaly detection via hierarchical memory networks." In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 201-218. Cham: Springer Nature Switzerland, 2023.
- [CVPR'23] Wang, Haoyu, Guansong Pang, Peng Wang, Lei Zhang, Wei Wei, and Yanning Zhang. "Glocal Energy-based Learning for Few-Shot Open-Set Recognition." In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7507-7516. 2023.
- [DSAA'23] Li, Jiayi, Guansong Pang, Ling Chen, and Mohammad-Reza Namazi-Rad. "HRGCN: Heterogeneous Graph-level Anomaly Detection with Hierarchical Relation-augmented Graph Neural Networks." In *2023 IEEE 10th International Conference on Data Science and Advanced Analytics (DSAA)*, pp. 1-10. IEEE, 2023.
- [NeurIPS'23] Hezhe Qiao and Guansong Pang. "Truncated Affinity Maximization: One-class Homophily Modeling for Graph Anomaly Detection." In: *NeurIPS'23*.