



Precise Object Localization using eXplainable AI Methods

PhD Student, Cycle 38
Muhammad Rashid¹

Academic Supervisor:
Elvio G. Amparore¹

Industrial Supervisors:
Enrico Ferrari²
Damiano Verda²


eXplainable AI (XAI)

Introduction

- Why and What is eXplainable Artificial Intelligence(XAI)
- Semantic Analysis of Machine Learning Models
- Explanations from Image Classification Models
- Explanations from Anomaly Detection Systems
- Class Activation Map based methods
- Better evaluation of AI systems using XAI
- Precise Object Localization in Multiple Application

Research Questions & Solutions


- Under representation of synthetic data (RQ1)
- Anomaly map finds precise Localization of Anomaly (RQ2)
- Shapely Values generation care about data? (RQ3)
- Anomaly Score is best choice for decision? (RQ4)
- SoTA evaluation metrics perform faithful evaluation? (RQ5)



RQ1: Under Representation of Synthetic Neighborhood

(a) Under Representation of Synthetic Neighborhood

- Perturbation based methods use synthetic neighborhood.
- Does SoTA XAI methods represent 100% of synthetic neighborhood?
- Under-representation causes confused and meaningless explanations.




(c) Proposed

Masks Generation mask vector x are sampled $x[i] \sim B(0.5), 1 \leq i \leq k$

Dependent Variables $Y = \{f(\xi_x) \mid \xi_x \in N(\xi)\}$

Distance Function ξ should weight more


$$w_x = \exp\left(-\frac{D(x)}{\sigma^2}\right)$$




Stratified sampling

- Partitioning stratified on the number of masked superpixels.
- $\mathcal{X}^{(i)}$ = set of possible masks having $|x| = i$
- Stratum i size is known a priori: $|\mathcal{X}^{(i)}| = \binom{k}{i}, 0 \leq i \leq k$



(d) Results



Monte Carlo sampling (default)



Stratified sampling



(e) Experimental Evaluation

$$\bullet CV(\beta) = \frac{\sigma_\beta}{\mu_\beta} \bullet RC(Y) = \frac{IQR_{1-99}(Y)}{f(\xi)}$$

Conclusions

- Reformulation of LIME sampling strategy.
- Practical improvements, at no additional cost.

The 38th Annual AAAI Conference on Artificial Intelligence Using Stratified Sampling to Improve LIME Image Explanations

Muhammad Rashid¹, Elvio G. Amparore¹, Enrico Ferrari², Damiano Verda²

Source Code

• pip install lime stratified
• github/rashidrao-pk/lime_stratified




RQ2: Anomaly map finds precise Localization of Anomaly? Integration of XAI for Anomalies in Images

(a) Problem Statement


Separating real anomaly and background noise (if any)

- Case study Explaining Anomaly Detection(AD) Systems
- Verifying anomaly detection systems
- Precise anomaly localization using XAI
- Use of One Class Classification based Self Supervised Learning for Anomaly Detection

(c) Detecting Anomalies



(d) Anomaly Localization – Role of XAI



(f) Conclusions

- XAI methods are relevant in finding the true drivers behind AI systems using techniques like classification and/or anomaly detection.
- Case study based on reconstruction error maps generated from VAE-GAN models.
- Multiple XAI techniques to separate the reconstruction error (noise) from the anomaly (if any).
- A sample may be detected as anomalous for the wrong reasons, yet this misbehaviour may not be detectable from the information provided by the anomaly detection system alone → Role of XAI!

(b) Training of Variational AutoEncoder Generative AI (VAE – GAN)

Encoding

- $z = e(\xi)$

Anomaly Map

- $m = |gs(\xi) - gs(\xi')|$

Decoding

- $\xi' = d(z)$

Non-defective object

Input image ξ


VAE-GAN model

Reconstructed image ξ'


Anomaly map

$$\text{Global Threshold} \quad \tau = \arg\max \sqrt{\text{TPR}(\tau) \times (1 - \text{FPR}(\tau))}$$

True Positive Rate : Anomalous as anomalous
False Positive Rate : Normal as anomalous



(e) Results



*SHAP: Scott, M. A unified approach to interpreting model predictions.

*LIME: Marco Túlio. 'Why should I trust you?' Explaining the predictions of any classifier.'

2nd World Conference on eXplainable Artificial Intelligence

Can I Trust my Anomaly Detection system? A case study based on eXplainable AI

Muhammad Rashid¹, Elvio G. Amparore¹, Enrico Ferrari², Damiano Verda²

Source Code

github/rashidrao-pk/anomaly_detection_trust_case_study



(a) Precise Object Localization

- Does SoTA XAI methods provide precise object localization?

Binary Hierarchy of Owen Approximation of Shapley Values


- Use of Data Aware approach - Binary Partition Tree (BPT)

- Feedback loop between Segmentation and Shapely Values Generation


$$\Delta_i(S) = \nu(S \cup \{i\}) - \nu(S)$$

(b) How Binary Partition Tree(BPT) works on Image?


$$dist(T_i, T_j) = dist_{color}^2(T_i, T_j) \cdot area(T_i, T_j) \cdot \sqrt{perim(T_i, T_j)}$$




(c) Comparison of BPT & Axis Aligned(AA) Hierarchies



(d) Evaluation Methodology



(e) Broad applicability of ShapBPT



(f) Performance Measures

$$AUC^+ = \int_0^1 \nu(S^{[q]}) dq$$

$$AUC^- = \int_0^1 \nu(N \setminus S^{[q]}) dq$$

$$MSE^+ = \int_0^1 (\nu(S^{[q]}) - \eta(S^{[q]}))^2 dq,$$

$$MSE^- = \int_0^1 (\nu(N \setminus S^{[q]}) - \eta(N \setminus S^{[q]}))^2 dq$$

$$AUOI = \int_0^1 J(S^{[q]}, G) dq$$

$$maxIoU = \max_{q \in [0,1]} (J(S^{[q]}, G))$$

(g) Experimental Evaluation

- Performance Curve based evaluation

GroundTruth based evaluation

• Controlled setup for exact IoU

- Multiple replacement values

• ImagNet Dataset

• ResNet50, VGG16, Swin-ViT

- Multiple Applications

• Object Localization

• Facial Attributes Localization

• Anomaly Localization

	AA	BPT	Ground truth	BPT
b=10				
b=100				

IoU=0.374 IoU=0.762

Research Trainings & Further Plans

Training Activities

- PhD Course Work (Completed)
- Conferences (ECML, AAAI-24, XAI World)
- External Courses (Oxford ML School, Google)
- Teaching Collaboration (Cyber-Security Course)
- Served as Reviewer in A & B ranked conferences.

Enterprise Time

First Half:

- Basic Training Courses
- Review, Learn and Implement
 - One Class Classification
 - AD using Hand Crafted-Features
 - AD using Self Supervised Learning
- Example/Prototype based explanations.
- Improvement of Anomaly Detection systems.
- Improved Anomaly Score Functions (RQ4).