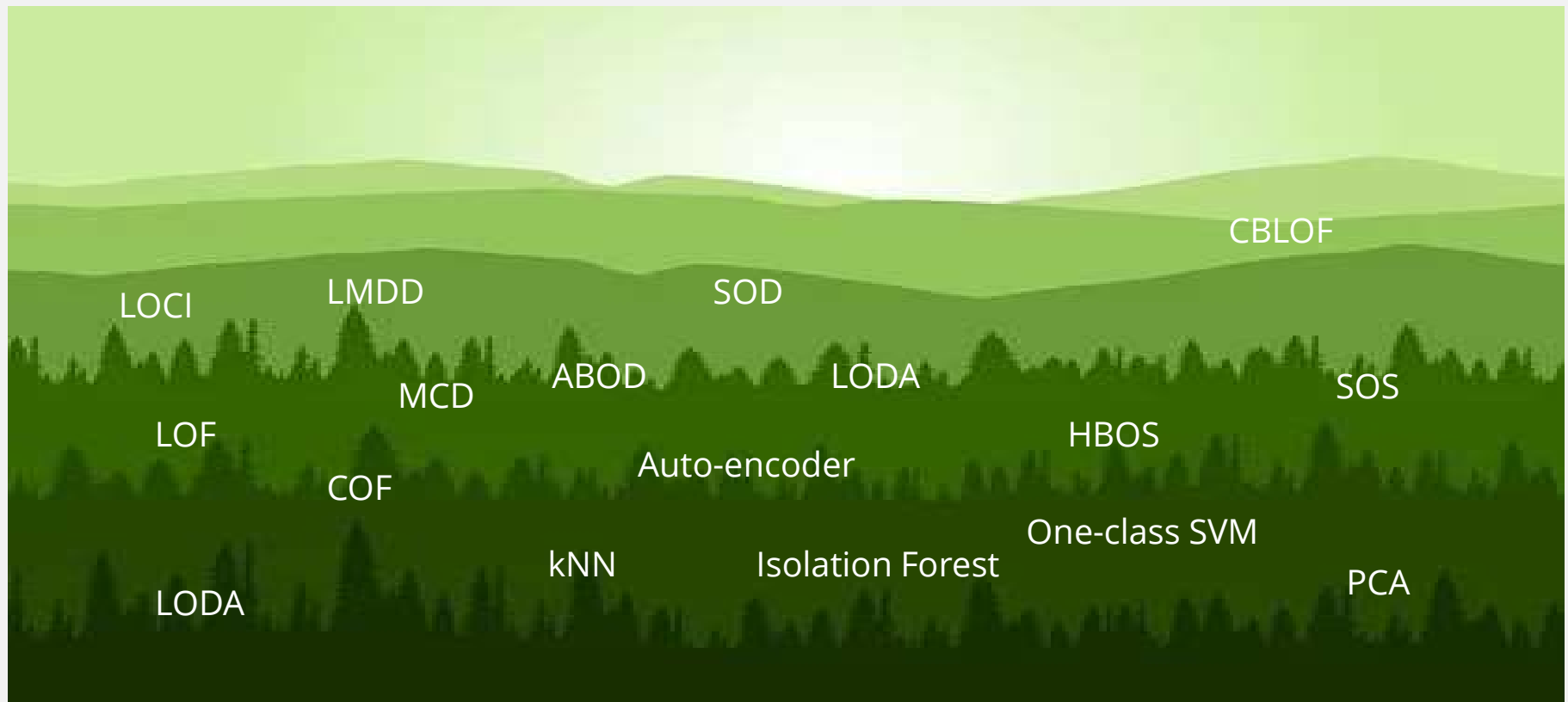# Comparing unsupervised anomaly detection algorithms

Roel Bouman

# ANOMALY DETECTION IN PREDICTIVE MAINTENANCE

- One of the first steps in setting up predictive maintenance is the detection of failures in historical data

- Historical data is in practice often of low quality
    - Logs are incomplete
    - Timestamps are off by minutes, hours or days
    - No good "labelling" of data exist (exact moments of failure)
- Additionally, not everything you'd want to detect is present in historical data (rare events)

- Often, we have to resort to Anomaly Detection to detect faults or failures in the absence of labels

- Knowing exactly when faults or failures happen allows for further investigation/modelling
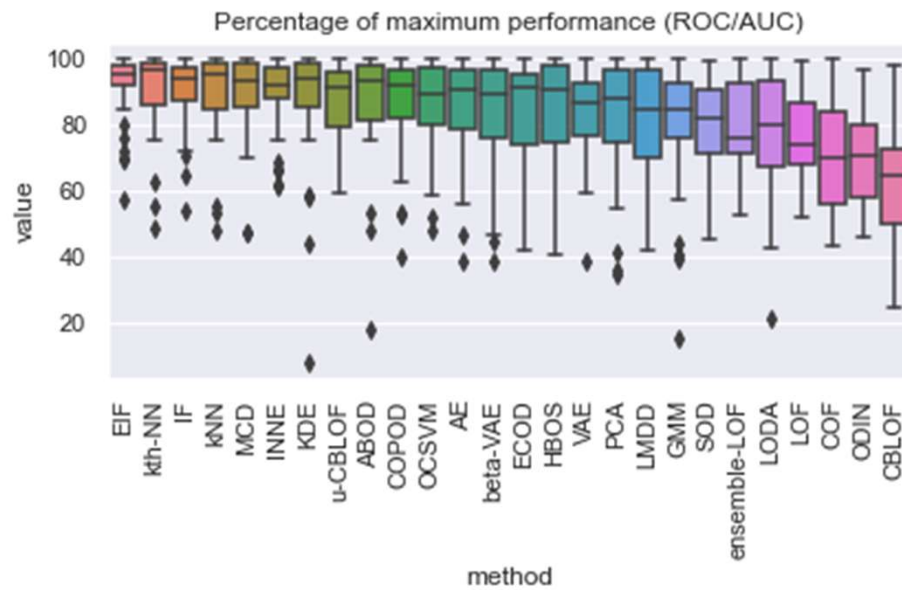
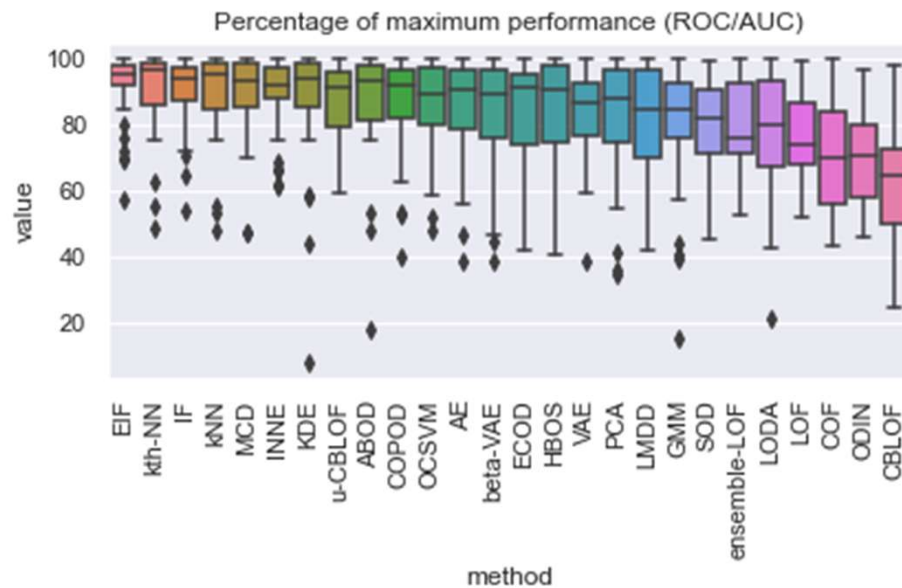# ANOMALY DETECTION ALGORITHMS: SEEING THE FOREST FOR THE TREES

# COMPARING UNSUPERVISED ANOMALY DETECTION ALGORITHMS

- For supervised classification, large scale comparison studies have been performed (Fernández-Delgado and Amorim)

- Yet, for unsupervised anomaly detection, fairly little comparitive research has been done
  - Goldstein and Uchida 2016 (19 algorithms, 10 datasets, no statistical comparison)
  - Campos et al. 2016 (12 algorithms, 11 datasets)

- Our study: 26 algorithms on 38 real-world tabular datasets (currently)
  - Use Imam-Davenport to check for presence of significant differences
  - Nemenyi-Friedman for pairwise testing
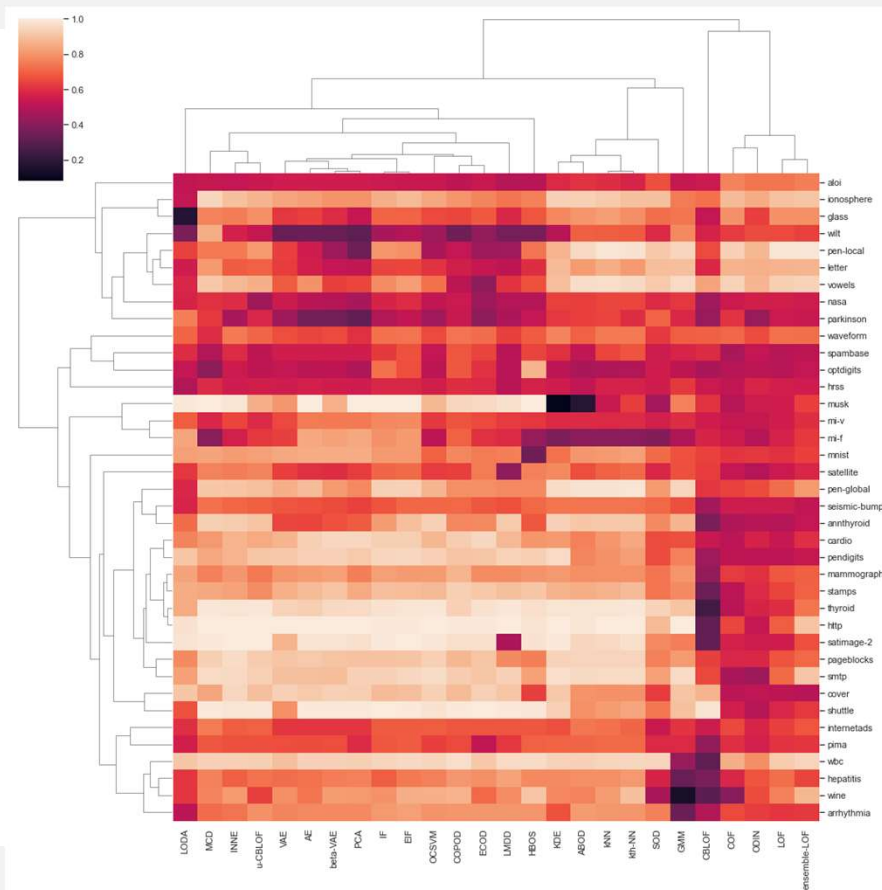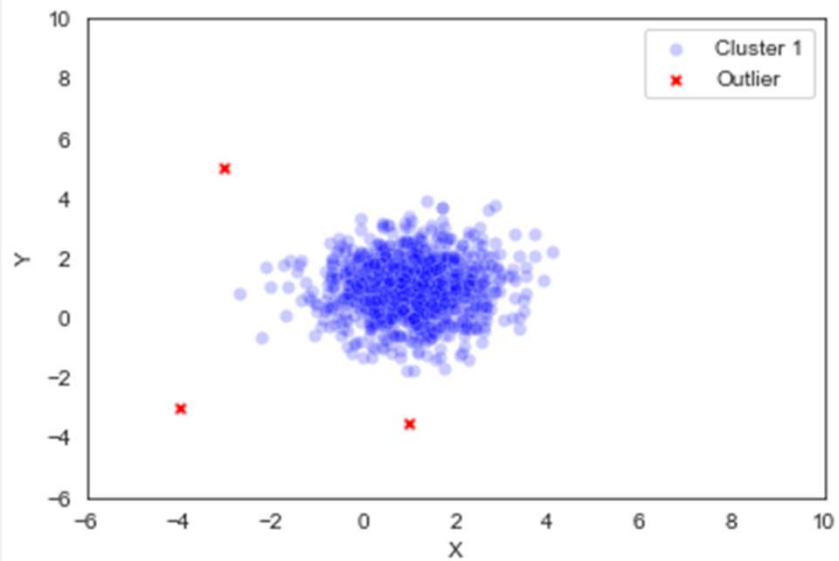
# COMPARING OVERALL PERFORMANCE



Percentage of maximum performance (ROC/AUC)

# COMPARING OVERALL PERFORMANCE



Percentage of maximum performance (ROC/AUC)

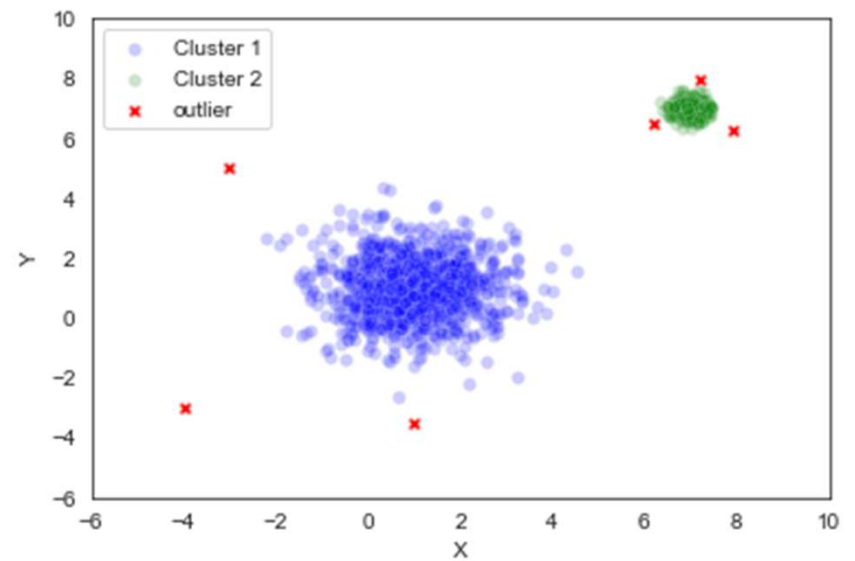| | CBLOF | ODIN | COF | LOF | LODA | ensemble-LOF | SOD | LMDD | PCA | Mean AUC |
|---|---|---|---|---|---|---|---|---|---|---|
| EIF | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | 0.798 |
| kth-NN | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | 0.787 |
| IF | ++ | ++ | ++ | ++ | ++ | | ++ | | | 0.787 |
| kNN | ++ | ++ | ++ | ++ | ++ | ++ | ++ | + | + | 0.782 |
| MCD | ++ | ++ | ++ | ++ | + | | | | | 0.779 |
| INNE | ++ | ++ | ++ | ++ | + | | | | | 0.778 |
| KDE | ++ | ++ | ++ | ++ | ++ | | ++ | | | 0.764 |
| u-CBLOF | ++ | ++ | ++ | | | | | | | 0.763 |
| ABOD | ++ | ++ | ++ | + | | | | | | 0.758 |
| COPOD | ++ | ++ | ++ | | | | | | | 0.754 |
| OCSVM | ++ | ++ | ++ | | | | | | | 0.745 |
| AE | ++ | ++ | | | | | | | | 0.744 |
| beta-VAE | ++ | | | | | | | | | 0.732 |
| ECOD | ++ | + | | | | | | | | 0.732 |
| HBOS | ++ | + | | | | | | | | 0.731 |
| VAE | ++ | | | | | | | | | 0.725 |
| PCA | + | | | | | | | | | 0.723 |
| LMDD | + | | | | | | | | | 0.706 |
| GMM | ++ | | | | | | | | | 0.698 |
| SOD | | | | | | | | | | 0.694 |
| ensemble-LOF | + | | | | | | | | | 0.688 |
| LODA | | | | | | | | | | 0.684 |
| LOF | | | | | | | | | | 0.652 |
| COF | | | | | | | | | | 0.612 |
| ODIN | | | | | | | | | | 0.605 |
| CBLOF | | | | | | | - | | - | - | 0.532 |

# TWO-WAY CLUSTERING OF ALGORITHMS AND DATASETS

## PROPERTIES OF ANOMALIES: LOCAL AND GLOBAL DENSITY ANOMALIES

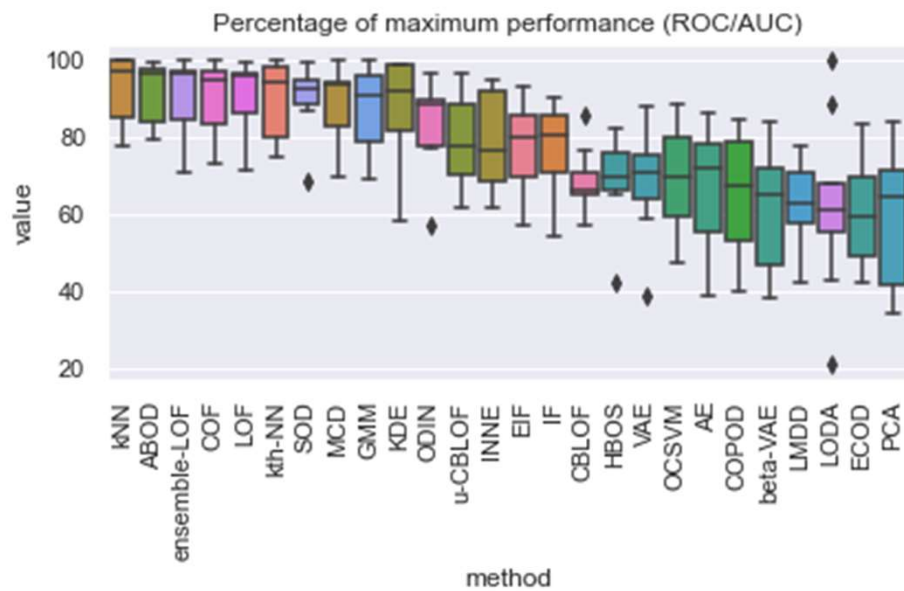# REPEATING THE COMPARISON FOR GLOBAL PROBLEMS (29 DATASETS)



Percentage of maximum performance (ROC/AUC)

| | CBLOF | COF | ODIN | LOF | ensemble-LOF | SOD | GMM | LODA | Mean AUC |
|---|---|---|---|---|---|---|---|---|---|
| EIF | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | 0.844 |
| IF | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | 0.833 |
| COPOD | ++ | ++ | ++ | ++ | ++ | ++ | | | 0.819 |
| INNE | ++ | ++ | ++ | ++ | ++ | ++ | | | 0.811 |
| AE | ++ | ++ | ++ | ++ | ++ | ++ | | | 0.804 |
| beta-VAE | ++ | ++ | ++ | ++ | + | ++ | | | 0.802 |
| ECOD | ++ | ++ | ++ | ++ | ++ | ++ | | | 0.801 |
| OCSVM | ++ | ++ | ++ | ++ | ++ | ++ | | | 0.801 |
| PCA | ++ | ++ | ++ | ++ | | | | | 0.800 |
| kth-NN | ++ | ++ | ++ | ++ | ++ | ++ | | | 0.794 |
| u-CBLOF | ++ | ++ | ++ | ++ | + | ++ | | | 0.793 |
| MCD | ++ | ++ | ++ | ++ | ++ | ++ | | | 0.792 |
| kNN | ++ | ++ | ++ | ++ | ++ | ++ | | | 0.781 |
| HBOS | ++ | ++ | ++ | ++ | + | ++ | | | 0.778 |
| VAE | ++ | ++ | ++ | ++ | | | | | 0.772 |
| KDE | ++ | ++ | ++ | ++ | ++ | ++ | | | 0.770 |
| LMDD | ++ | ++ | ++ | ++ | | | | | 0.767 |
| ABOD | ++ | ++ | ++ | ++ | | | | | 0.753 |
| LODA | | | | | | | | | 0.736 |
| GMM | | + | | | | | | | 0.680 |
| SOD | | | | | | | | | 0.673 |
| ensemble-LOF | | | | | | | | | 0.659 |
| LOF | | | | | | | | | 0.610 |
| ODIN | | | | | | | | | 0.566 |
| COF | | | | | | | | - | 0.556 |
| CBLOF | | | | | | | | | 0.508 |

Percentage of maximum performance (ROC/AUC)

| | PCA | LMDD | beta-VAE | ECOD | LODA | COPOD | AE | OCSVM | VAE | HBOS | CBLOF | Mean AUC |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| kNN | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | 0.787 |
| ABOD | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | ++ | 0.773 |
| ensemble-LOF | ++ | ++ | ++ | ++ | | | + | | | + | + | 0.770 |
| LOF | ++ | ++ | ++ | ++ | | + | ++ | | | + | + | 0.768 |
| kth-NN | ++ | ++ | ++ | ++ | + | + | ++ | | | ++ | + | 0.767 |
| COF | ++ | ++ | ++ | ++ | | | | | | | | 0.767 |
| SOD | ++ | ++ | ++ | | | | | | | | | 0.751 |
| GMM | ++ | ++ | + | | | | | | | | | 0.748 |
| KDE | ++ | ++ | ++ | ++ | | + | ++ | | | + | + | 0.746 |
| MCD | + | + | | | | | | | | | | 0.740 |
| ODIN | | | | | | | | | | | | 0.714 |
| INNE | | | | | | | | | | | | 0.685 |
| u-CBLOF | | | | | | | | | | | | 0.682 |
| EIF | | | | | | | | | | | | 0.671 |
| IF | | | | | | | | | | | | 0.658 |
| CBLOF | | | | | | | | | | | | 0.601 |
| HBOS | | | | | | | | | | | | 0.600 |
| VAE | | | | | | | | | | | | 0.594 |
| OCSVM | | | | | | | | | | | | 0.590 |
| AE | | | | | | | | | | | | 0.574 |
| COPOD | | | | | | | | | | | | 0.571 |
| LODA | | | | | | | | | | | | 0.538 |
| ECOD | | | | | | | | | | | | 0.537 |
| beta-VAE | | | | | | | | | | | | 0.537 |
| LMDD | | | | | | | | | | | | 0.535 |
| PCA | | | | | | | | | | | | 0.507 |

# CONCLUSION

- We've found a subset of algorithms which work well on various types of anomalies:

  - kNN works well on the entire collection of datasets, as well as on both local and global anomalies

  - Extended Isolation Forest works best on global anomalies

  - KNN works best on local anomalies

- The current benchmark datasets require more analysis to study which properties the contained anomalies have

Radboud University

# FUTURE WORK

- Extending the benchmark and keeping it up-to-date

- There are no tests to see what properties the anomalies within a certain dataset have

- Look further into different properties of algorithms:

  - Multidimensional vs. Unidimensional

  - Enclosed vs. Peripheral

  - Isolated vs. Clustered

- Look into hyperparameter/initialisation stability

# QUESTIONS