

A Transductive Forest for Anomaly Detection with Few Labels

<https://github.com/jzha968/transForest>

Jingrui Zhang, Ninh Pham, Gillian Dobbie

University of Auckland



ECML-PKDD, Sep 19, 2023

Anomaly detection

- Definition:

- Identify rare items, events, or observations which deviate significantly from the majority of the data

- Unsupervised methods:

- Common assumption: “Anomalies are on *sparse* regions while normal points are on *dense* regions.”

- Challenges in high dimensions:

- The common unsupervised prior is not true due to irrelevant features
- There are exponential number of subspaces to investigate

This work

- Research question: How could we find relevant subspaces to identify anomalies given few labels?
- Solution: Semi-supervised TransForest
 - Push the classification boundaries towards sensitive subspaces containing both normal and abnormal points
 - Provide feature importance ranking
 - Competitive with recent semi-supervised SOTA with 2% labeled data

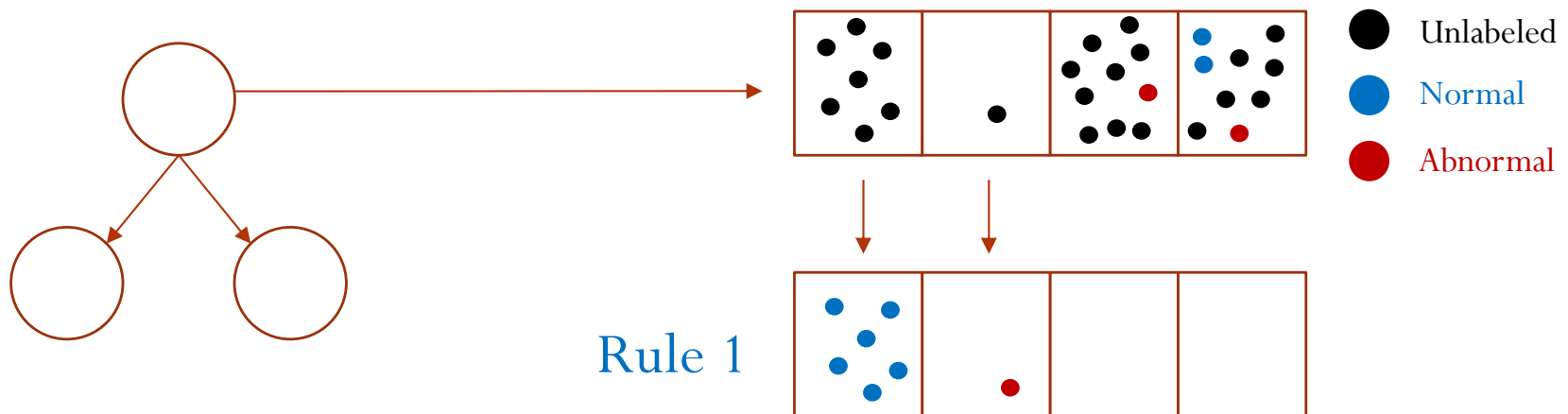
How TransForest works?

- Setting:
 - Collection of trees (similar to iForest, Random Forest, Extra Trees)
 - Tree's input: **all** available labeled points + a subset of unlabeled points
- Training: Pseudo-labeling + Extra Trees learning
 - For a randomly selected feature \mathbf{f}_i , build a **histogram**, and use label information to **pseudo-label** unlabeled points in each bin
 - For a random feature \mathbf{f}_i and a **random** cut \mathbf{v}_i , compute the InfoGain based on the labeled and pseudo-labeled points; then, select $(\mathbf{f}_i, \mathbf{v}_i)$ that maximizes InfoGain

Pseudo-label points in a bin \mathbf{B}

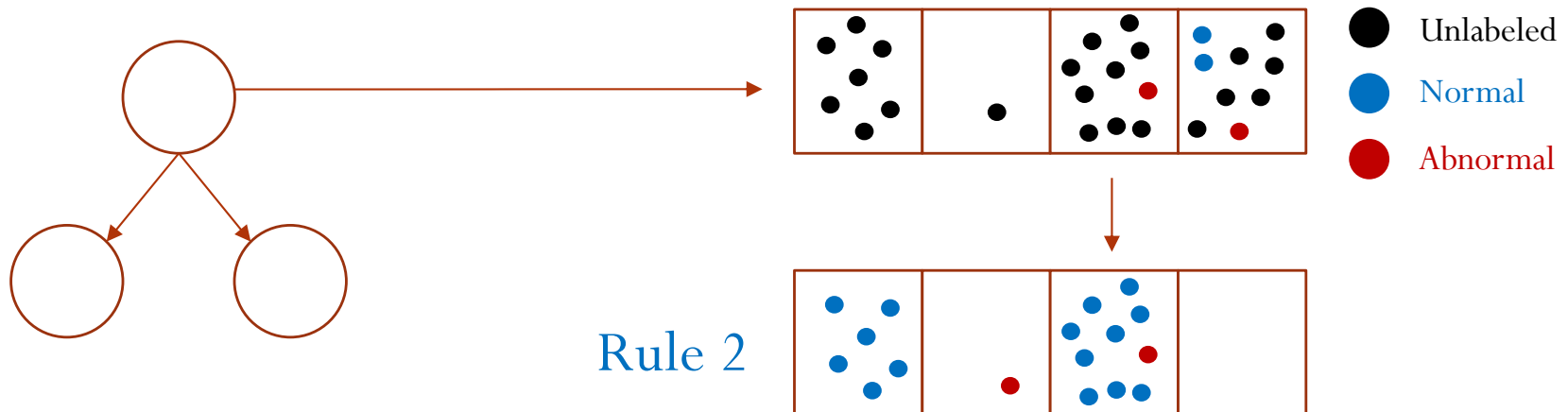
- Rule 1 (unsupervised):

- If \mathbf{B} has no labeled points, resort unsupervised prior given a threshold Δ
 - If \mathbf{B} is **dense** ($|\mathbf{B}| \geq \Delta$), all points in \mathbf{B} are **normal**
 - If \mathbf{B} is **sparse** ($|\mathbf{B}| < \Delta$), all points in \mathbf{B} are **abnormal**



Pseudo-label points in a bin \mathbf{B}

- Rule 2 (anomalies are rare):
 - If \mathbf{B} is dense and contains only labeled anomalies
 - \mathbf{B} will have $0.9 |\mathbf{B}|$ normal and $0.1 |\mathbf{B}|$ abnormal points

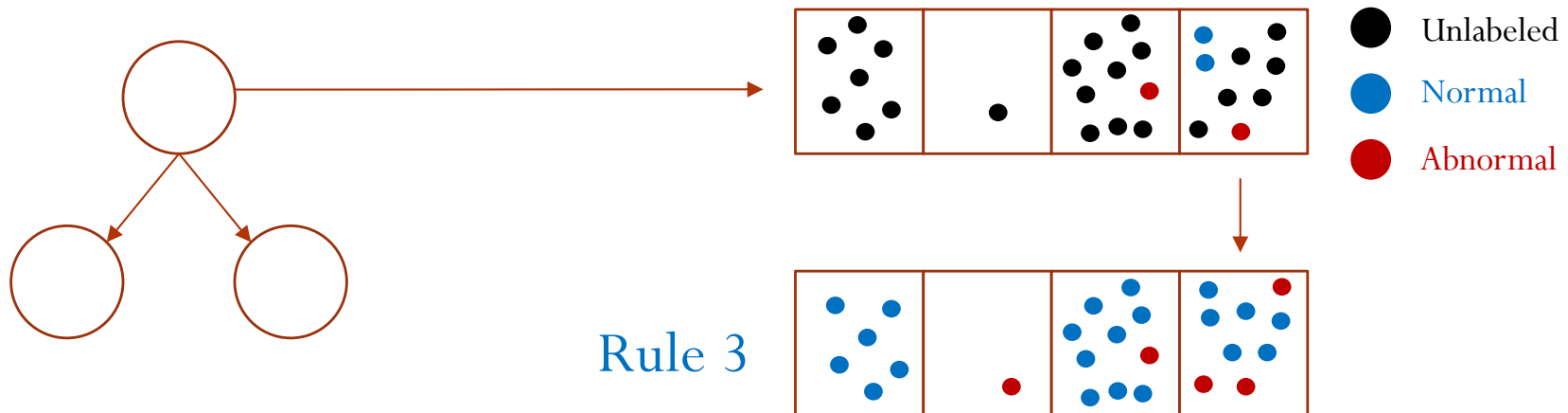


Pseudo-label points in a bin \mathbf{B}

- Rule 3 (semi-supervised):

- If \mathbf{B} has \mathbf{m}_0 normal points and \mathbf{m}_1 anomalies

- \mathbf{B} will have $\frac{\mathbf{m}_0}{\mathbf{m}_0 + \mathbf{m}_1} |\mathbf{B}|$ normal and $\frac{\mathbf{m}_1}{\mathbf{m}_0 + \mathbf{m}_1} |\mathbf{B}|$ abnormal points



Pseudo-label points in a bin \mathbf{B}

- Rule 1 (unsupervised):

- If \mathbf{B} has no labeled points, resort unsupervised prior given a threshold Δ
 - If \mathbf{B} is **dense** ($|\mathbf{B}| \geq \Delta$), all points in \mathbf{B} are **normal**
 - If \mathbf{B} is **sparse** ($|\mathbf{B}| < \Delta$), all points in \mathbf{B} are **abnormal**

hyperparameters

- Rule 2 (anomalies are rare):

- If \mathbf{B} is dense and contains only labeled anomalies
 - \mathbf{B} will have $0.9 |\mathbf{B}|$ normal and $0.1 |\mathbf{B}|$ abnormal points

hyperparameters

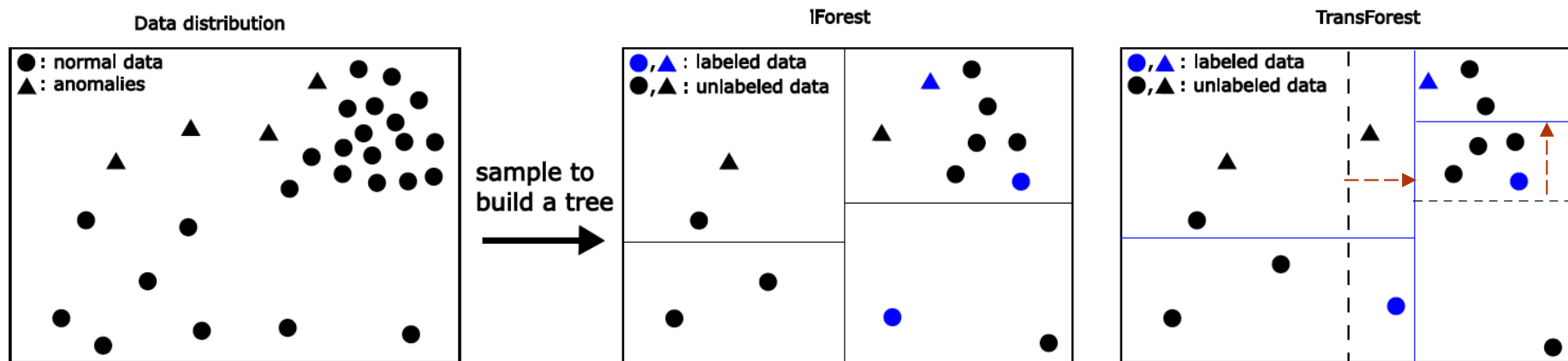
- Rule 3 (semi-supervised):

- If \mathbf{B} has \mathbf{m}_0 normal points and \mathbf{m}_1 anomalies
 - \mathbf{B} will have $\frac{\mathbf{m}_0}{\mathbf{m}_0 + \mathbf{m}_1} |\mathbf{B}|$ normal and $\frac{\mathbf{m}_1}{\mathbf{m}_0 + \mathbf{m}_1} |\mathbf{B}|$ abnormal points

TransForest vs iForest

- TransForest vs iForest:

- TransForest pushes boundaries towards sensitive areas



- Testing:

- Using **pathLength** as anomaly score (similar to iForest) with **adjustment**
 - Labeled anomalies on isolated node: **pathLength** = 1
 - Labeled normal points on isolated node: **pathLength** = $\log(s)$

Hyperparameter setting

- Hyperparameters:

- Density threshold of node S : $\Delta = 0.1 |S|$
- Number of histogram bins of node S : $\log(|S|) + 1$
- Number of trees: $t = 100$
- Number of trials to find the best split: $k = 10$
- Tree size: $s = \max(256, 2 n_l)$ where $n_l = \# \text{ labels}$

- Complexity:

- Training in $O(t s k \log(s))$ and testing in $O(n t \log(s))$ as similar as iForest

- Feature importance ranking:

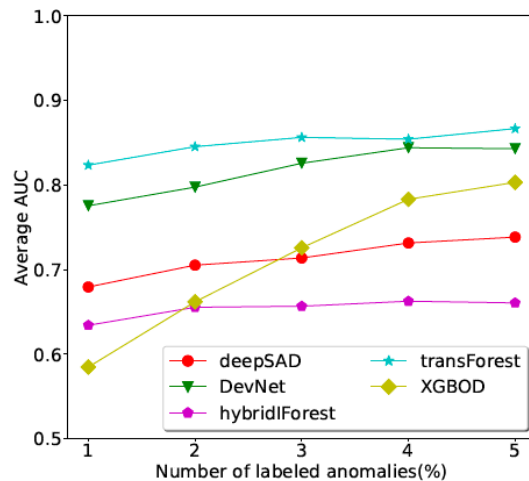
- Derived from the estimated InfoGain

Experiment

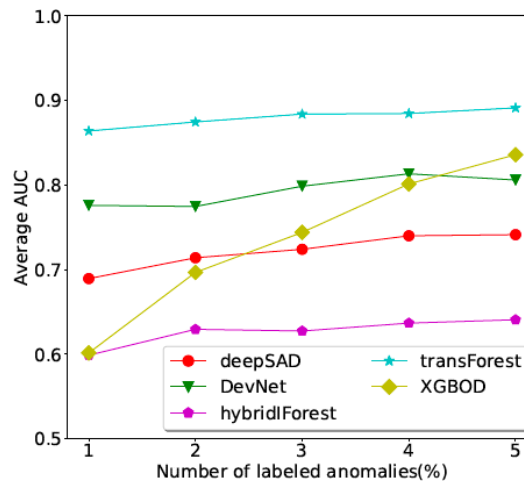
- Datasets:
 - 15 tabular and 20 continuous real-world data sets
- Experiment setting:
 - Randomly take 2% - 10% labeled data for training (# labeled normal points = # labeled anomalies) and use the rest for testing
 - Measure average AUC for 10 runs
- Semi-supervised competitors:
 - Tree and boosting: Hybrid iForest, XGBOD
 - Deep learning: DeepSAD, DevNet

Experiment

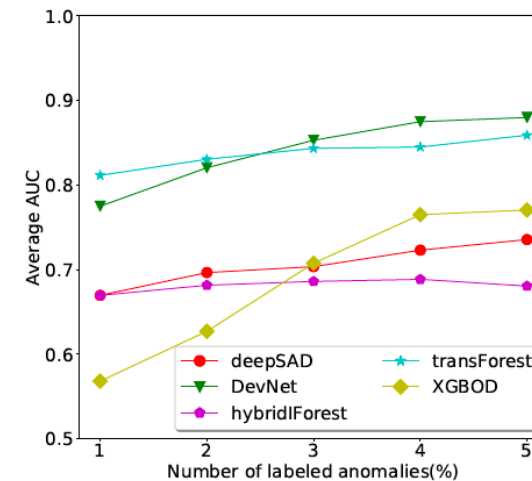
- Avg AUC using 2% - 10% labels on various data sets:



(a) All datasets.

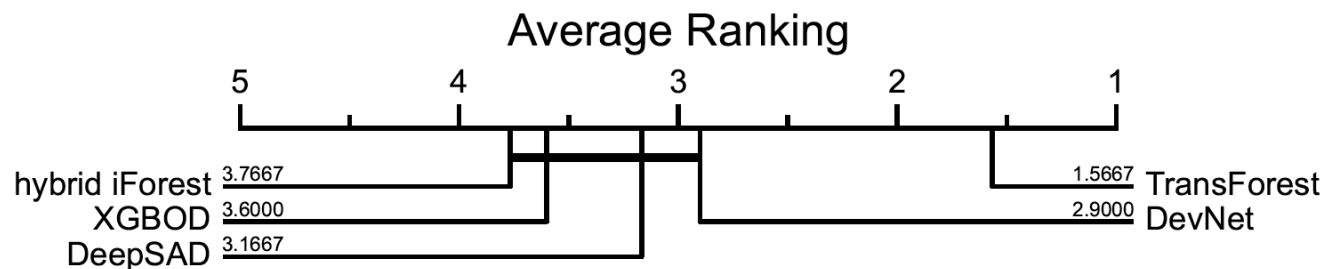


(b) Tabular datasets.



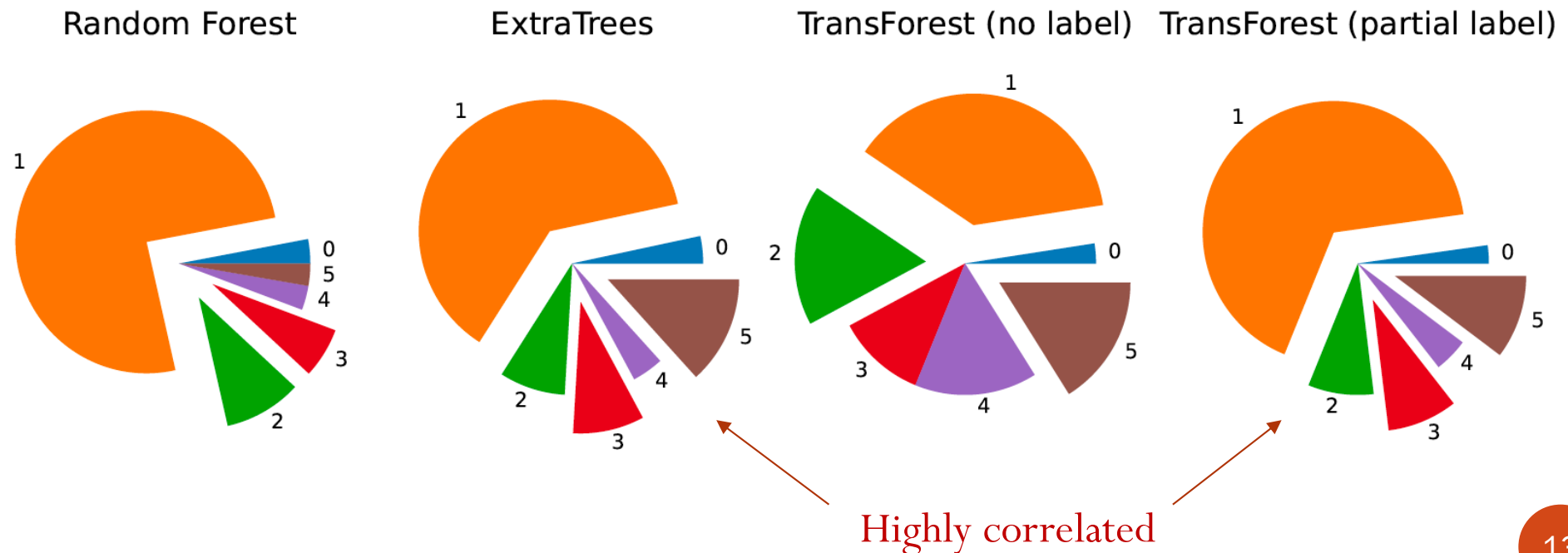
(c) CV & NLP datasets.

- Critical diagram using 2% labels on all data sets:



Feature importance ranking

- Annnthyroid ($n = 7200$, $d = 6$, # anomalies = 534):
 - We use $\sim 2\%$ of labeled points (10 anomalies and 10 normal points)
 - Number is dimension index
 - Wedge size reflects feature importance



Conclusion

- Few labels are useful for identifying relevant features for anomaly detection
- TransForest:
 - Simple but competitive with other semi-supervised models
 - Consistent importance feature ranking with supervised models on low-dimensional data sets
 - Robust against irrelevant features
 - <https://github.com/jzha968/transForest>