Anomaly detection and explanation

Martin Kopp

Czech Technical University in Prague

Cisco Systems, Cognitive Research Team in Prague

Institute of Computer Science, Academy of Sciences of the Czech Republic

March 12, 2015

《曰》 《聞》 《臣》 《臣》 三臣 …

Outline

- Anomaly detection
- Anomaly explanation
 - Sapling random forests
 - minimal explanation
 - maximal explanation
 - rules aggregation
- Clustering
 - voting vectors
 - feature deviations
 - evaluation



Outline

Anomaly detection

- Anomaly explanation
 - Sapling random forests
 - minimal explanation
 - maximal explanation
 - rules aggregation
- Clustering
 - voting vectors
 - feature deviations
 - evaluation



Anomaly detection is about ...





(a)

... point of view.





÷.

ヘロマ ヘヨマ ヘロマ

Anomaly explanation

Clustering

Anomaly detection

Anomaly in crowd







¹www.svcl.ucsd.edu/projects/anomaly/

◆□ → ◆□ → ◆ □ → ◆ □ → ○ □

Network security

- $\bullet\,$ typical proportion of anomalies is 1-0.1%
- 0.5 million data points \rightarrow 1000 anomalies

Particle physics

- typical proportion of anomalies is $10^{-3} 10^{-4}\%$
- $\bullet~2$ million data points \rightarrow 100 anomalies



Network security

- typical proportion of anomalies is 1-0.1%
- 0.5 million data points \rightarrow 1000 anomalies

Particle physics

- typical proportion of anomalies is $10^{-3} 10^{-4}\%$
- 2 million data points \rightarrow 100 anomalies



Anomaly detection problem statement

" An outlier is an observation which deviates so much from the other observations as to arouse suspicions that it was generated by a different mechanism."



¹Hawkins 1980 - Identification of outliers

◆□ → ◆□ → ◆ □ → ◆ □ → ○ □

Anomaly detection problem statement

- Defining a normal region for every possible normal behaviour is very difficult.
- The boundary between normal and anomalous behaviour is often not precise.
- Some anomalous events often adapt to appear normally.
- Even normal behaviour may evolve over time.
- Obtaining labelled data for training and validation of models is usually a major issue.
- Often the data contains noise that tends to be similar to the actual anomalies and hence is difficult to distinguish and remove.



Anomaly detectors

- Statistical
- Linear
- Proximity based
 - cluster
 - distance
 - density
- odomain specific



Statistical anomaly detectors





Linear model based detectors





・ロット (雪) (日) (日)

Cluster based detectors





æ

・ロト ・ 四ト ・ ヨト ・ ヨト

Distance based detectors



https://baldscientist.wordpress.com/2013/02/02/is-free-will-a-matter-ofbeing-a-conscious-outlier/



э

Density based detectors





э

 $http://scikit-learn.org/stable/modules/outlier_detection.html < \texttt{m} \land \texttt{m} `$

Outline

Anomaly detection

- Anomaly explanation
 - Sapling random forests
 - minimal explanation
 - maximal explanation
 - Clustering
 - voting vectors
 - feature deviations
 - evaluation
- 4 Rules
 - voting vectors
 - feature deviations
 - evaluation



Anomaly explanation history

- Grubbs 1950 Anomaly detection ¹
- Knorr 1999 Question ²
- Dang 2013 Answer ³

¹Grubbs 1950 - *Sample criteria for testing outlying observations.* ²Knorr, Edwin M., and Raymond T. Ng. 1999 - *Finding intensional knowledge of distance-based outliers.*

AS

³Dang, Xuan Hong, et al. 2013 - *Local outlier detection with interpretation.*

Anomaly explanation

- Network security
 - attack vs. unscheduled backup
- Particle physics
 - Higgs boson vs. misconfiguration of equipment
- Astronomy
 - cosmic microwave background vs. pigeon nest
- Fraud detection
 - holiday vs. credit card fraud



Anomaly explanation

We have:

- dataset
- anomaly detection algorithm
- labelled suspicious samples

We want:

- examine the suspicious samples
- interpret them clearly
 - as a small subset of features.
 - as human readable set of rules



Anomaly explanation

We have:

- dataset
- anomaly detection algorithm
- labelled suspicious samples

We want:

- examine the suspicious samples
- interpret them clearly
 - as a small subset of features
 - as human readable set of rules





(a) In nature



Sapling Random Forest sapling



(b) In theoretical informatics



æ

ヘロト 人間 とくほ とくほ とう

- ensembles of specifically trained CARTs
- multiple trees per anomaly
- specifically made training sets -> grow sets
- trees are quite small -> saplings



Summary of the SRF for minimal explanation

Input: data $y \leftarrow anomalyDetector(data)$ for all data(y ==anomaly) do $G \leftarrow createGrowSet(size, method)$ $T \leftarrow trainTree(G)$ $SRF \leftarrow T$ end for extractRules(SRF)



Input: data





 $y \leftarrow anomalyDetector(data)$





Sapling Random Forest algorithm

 $G \leftarrow createGrowSet(size, method)$



Sapling Random Forest Grow set selection

A grow set \mathcal{G} contains an anomaly x^a and several normal samples $x^n \subseteq \mathcal{X}^n$.

- typical size $|\mathcal{G}| = 100$
- random selection
 - fast even in high dimensions
 - multiple trees can be grown -> robust
- k-nn selection
 - deterministic more trees are useless
 - slow in high dimensions
 - superior in low dimensions



Sapling Random Forest Grow set selection





・ロト ・ 理 ト ・ 理 ト ・ 理 ト

$T \leftarrow trainTree(G)$



Sapling Random Forest splitting criterion

Gini's index

$$\mathbf{G}_{\mathbf{i}} = 1 - p_a^2 - p_n^2,$$

Information gain

$$\arg \max_{h \in \mathcal{H}} \quad -\sum_{b \in \{L,R\}} \frac{|\mathcal{S}^b(h)|}{|\mathcal{S}|} H(\mathcal{S}^b(h)),$$



Sapling Random Forest splitting criterion

Simplified criterion

 $\arg\min_{h\in\mathcal{H}}|\mathcal{S}^a(h)|,$



Sapling Random Forest splitting criterion

Simplified criterion

 $\arg\min_{h\in\mathcal{H}}|\mathcal{S}^a(h)|,$

Maximal margin

 $\arg\max_{d\in\mathcal{D}}\max\min\mathcal{S}_d^n-x_d^a$

 $\arg\max_{d\in\mathcal{D}}\inf\mathcal{S}_d^n-x_d^a$



Sapling Random Forest

$SRF \leftarrow T$







Sapling Random Forest tree training





æ

・ロト ・聞 ト ・ ヨ ト ・ ヨ ト

Sapling Random Forest explaining an anomaly

extractRules(*SRF*)

 $C = x_2 > 2.2$





Sapling Random Forest explaining an anomaly

extractRules(SRF)

$$C = (x_2 > 2.2) \land (x_1 < -2.1)$$





Sapling Random Forest explaining an anomaly

extractRules(SRF) $C = (x_2 > 2.2) \land (x_1 < -2.1) \land (x_1 > 2.2) \land \dots$





Sapling Random Forest explaining an anomaly

The set of all possible rules is defined as $\mathcal{H} = \{h_{j,\theta} | j \in \{1, \dots, d\}, \theta \in \mathbb{R}\}$ where

$$h_{j,\theta}(x) = \begin{cases} +1 & \text{if } x_j > \theta \\ -1 & \text{otherwise} \end{cases}$$

- d ... number of features
- θ ... inner node threshold
- x_j...jth feature of sample x



Sapling Random Forest explaining an anomaly

The set of all possible rules is defined as $\mathcal{H} = \left\{h_{j,\theta} | j \in \{1, \dots, d\}, \theta \in \mathbb{R}\right\}$ where

$$h_{j, heta}(x) = egin{cases} +1 & ext{if } x_j > heta \ -1 & ext{otherwise} \end{cases}$$

Let $h_{j_1,\theta_1}, \ldots, h_{j_r,\theta_r}$ be the set of decisions taken in inner nodes on the path from the root to the leaf with the anomaly x^a . Then x^a is explained as conjunction of atomic conditions



Sapling Random Forest rules extraction

Rules in form:

$$C = (x_{j_1} > \theta_1) \land (x_{j_2} < \theta_2) \land \ldots \land (x_{j_t} > \theta_t)$$

We calculate groups sizes

$$r_{2j} = \sum_{C \in \mathcal{D}} \sum_{h \in C} I(j \in h, L)$$
$$r_{2j-1} = \sum_{C \in \mathcal{D}} \sum_{h \in C} I(j \in h, R)$$
$$I(j \in h) = \begin{cases} +1 & \text{if } < rule \\ -1 & \text{otherwise} \end{cases}$$



æ

ヘロト 人間 とくほ とくほとう

Sapling Random Forest rules extraction

and chose only k-most frequent, where

$$k = \arg\min_{k} \frac{1}{\sum_{j=1}^{2d} r_j} \sum_{j=1}^{k} r_j > \tau$$

Then we aggregate similar rules and chose the most strict thresholds.

$$h_j^R = \arg\min_{h\in\mathcal{H}_j^R}\theta_h$$



・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

Summary of the SRF for maximal explanation

```
y \leftarrow anomalyDetection(data)
for all data(y == anomaly) do
  f \leftarrow allFeatures
  while d < \tau do
     G \leftarrow createGrowSet(size, f)
     t \leftarrow trainTree(G)
     SRF \leftarrow SRF + t
     f \leftarrow f - topSplitFeature(t)
     D = nnDistance(G)
     d = D(anomaly)/max(D)
  end while
end for
extractRules(SRF)
```







æ



(a) minimal explanation



(b) maximal explanation

ヘロン ヘロン ヘビン ヘビン



æ





(b) minimal explanation

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト



æ

Anomaly explanation as feature selection





э

・ロト ・ 四ト ・ ヨト ・ ヨト

Anomaly explanation as feature selection





э

・ロト ・ 四ト ・ ヨト ・ ヨト

Outline

- Anomaly detection
- 2 Anomaly explanation
 - Sapling random forests
 - minimal explanation
 - maximal explanation
 - rules aggregation

Clustering

- voting vectors
- feature deviations
- evaluation



Clustering motivation

- Investigation of multiple anomalies at once
- Generalized anomaly groups
- Discovery of large scale anomalies
- Domain knowledge



Clustering Voting vectors

- binary vector
- tree voting
- TxA matrix
- sapling are anomaly specific
- sapling votes for similar anomalies



Clustering Voting vectors

Example of voting vectors





æ

Clustering Features deviation matrix

- deviation in feature ranges
- the most strict threshold is stored
- lower and upper boundary
- Tx2d matrix, but can be reduced



Clustering Voting

Example of features deviation matrix



Clustering results

Grow set size vs performance





æ

Clustering results

Number of clusters vs performance





æ

Conclusion and future work

Conclusion

- anomaly explanation
 - most important features
 - human readable rules
- arbitrary anomaly detector
- real time/data streams

Future work

- multi-dimensional anomalies
- cluster rules aggregation
- fuzzy rules



Thank you for your attention.

