

Data Mining

Munawar, PhD

12. Anomaly Detection





Outline

- Introduction
- Aspects of Anomaly Detection Problem
- Applications
- Different Types of Anomaly Detection

What are Anomalies?

- Anomaly is a pattern in the data that does not conform to the expected behavior
- Also referred to as outliers, exceptions, peculiarities, surprise, etc.
- Anomalies translate to significant (often critical) real life entities
 - Cyber intrusions
 - Credit card fraud

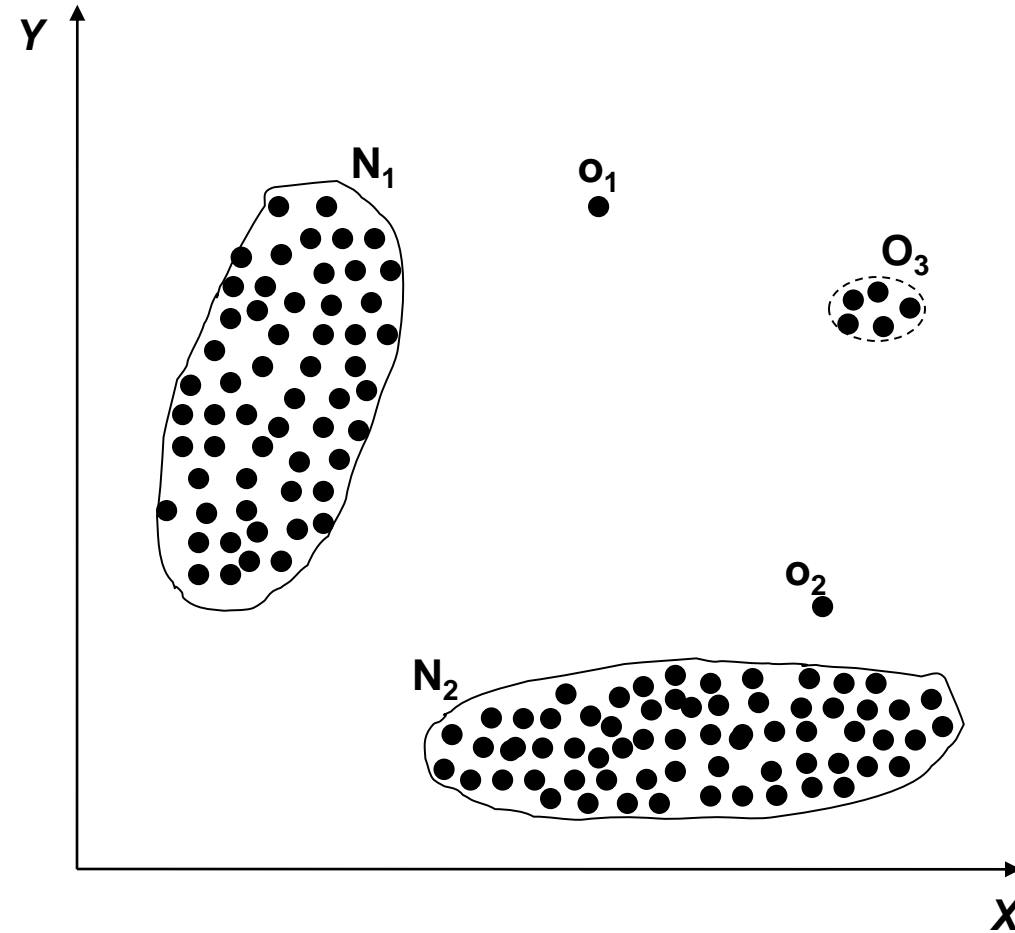
Real World Anomalies

- Credit Card Fraud
 - An abnormally high purchase made on a credit card
- Cyber Intrusions
 - A web server involved in *ftp* traffic



Simple Example

- N_1 and N_2 are regions of normal behavior
- Points o_1 and o_2 are anomalies
- Points in region O_3 are anomalies



Key Challenges

- Defining a representative normal region is challenging
- The boundary between normal and outlying behavior is often not precise
- The exact notion of an outlier is different for different application domains
- Availability of labeled data for training/validation
- Malicious adversaries
- Data might contain noise
- Normal behavior keeps evolving

Aspects of Anomaly Detection Problem

- Nature of input data
- Availability of supervision
- Type of anomaly: point, contextual, structural
- Output of anomaly detection
- Evaluation of anomaly detection techniques

Input Data

- Most common form of data handled by anomaly detection techniques is *Record Data*
 - Univariate
 - Multivariate

<i>Tid</i>	SrcIP	Start time	Dest IP	Dest Port	Number of bytes	Attack
1	206.135.38.95	11:07:20	160.94.179.223	139	192	No
2	206.163.37.95	11:13:56	160.94.179.219	139	195	No
3	206.163.37.95	11:14:29	160.94.179.217	139	180	No
4	206.163.37.95	11:14:30	160.94.179.255	139	199	No
5	206.163.37.95	11:14:32	160.94.179.254	139	19	Yes
6	206.163.37.95	11:14:35	160.94.179.253	139	177	No
7	206.163.37.95	11:14:36	160.94.179.252	139	172	No
8	206.163.37.95	11:14:38	160.94.179.251	139	285	Yes
9	206.163.37.95	11:14:41	160.94.179.250	139	195	No
10	206.163.37.95	11:14:44	160.94.179.249	139	163	Yes

Input Data – *Nature of Attributes*

- Nature of attributes
 - Binary
 - Categorical
 - Continuous
 - Hybrid

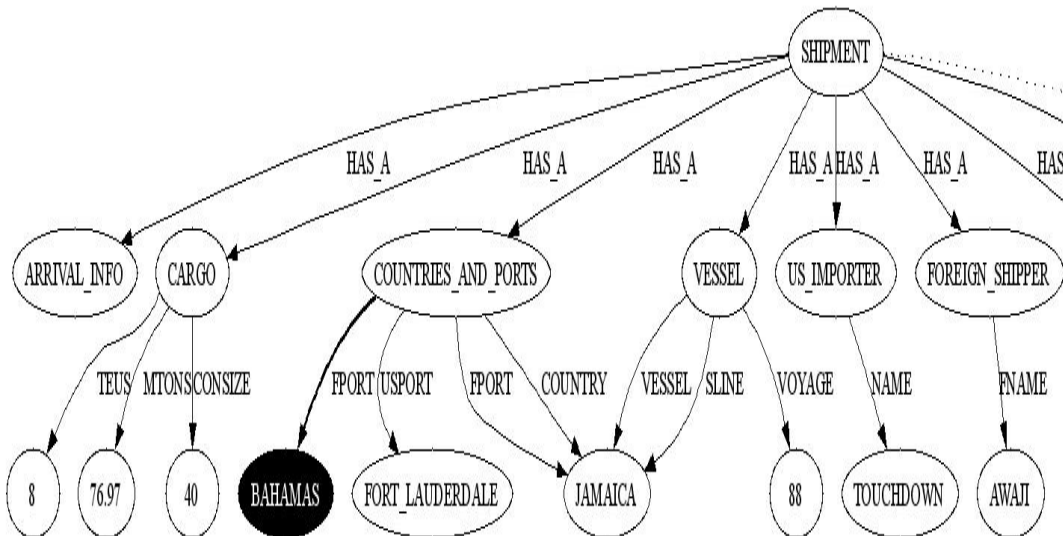
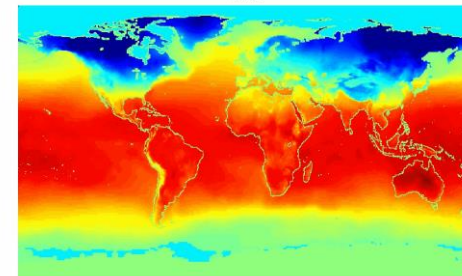
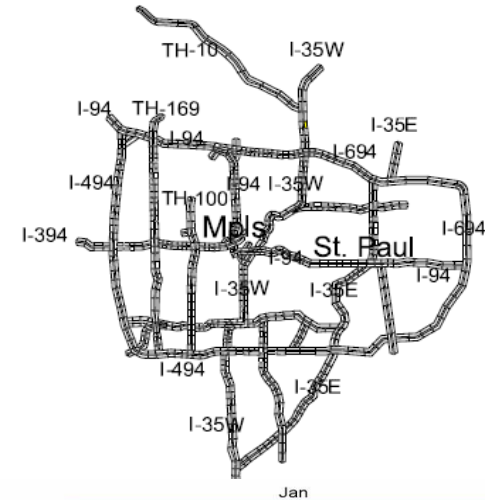
<i>Tid</i>	<i>SrcIP</i>	<i>Duration</i>	<i>Dest IP</i>	<i>Number of bytes</i>	<i>Internal</i>
1	206.163.37.81	0.10	160.94.179.208	150	No
2	206.163.37.99	0.27	160.94.179.235	208	No
3	160.94.123.45	1.23	160.94.179.221	195	Yes
4	206.163.37.37	112.03	160.94.179.253	199	No
5	206.163.37.41	0.32	160.94.179.244	181	No

categorical *continuous* *categorical* *continuous* *binary*

Input Data – *Complex Data Types*

- Relationship among data instances
 - Sequential
 - Temporal
 - Spatial
 - Spatio-temporal
 - Graph

```
GGTTCCGCCTTCAGCCCCGCGCC
CGCAGGGCCCGCCCCGCGCCGTC
GAGAAGGGCCCGCCTGGCGGGCG
GGGGGAGGCGGGGCGCCCGAGC
CCAACCGAGTCCGACCAGGTGCC
CCCTCTGCTCGGCCTAGACCTGA
GCTCATTAGGCGGCAGCGGACAG
GCCAAGTAGAACACGCGAAGCGC
```



Data Labels

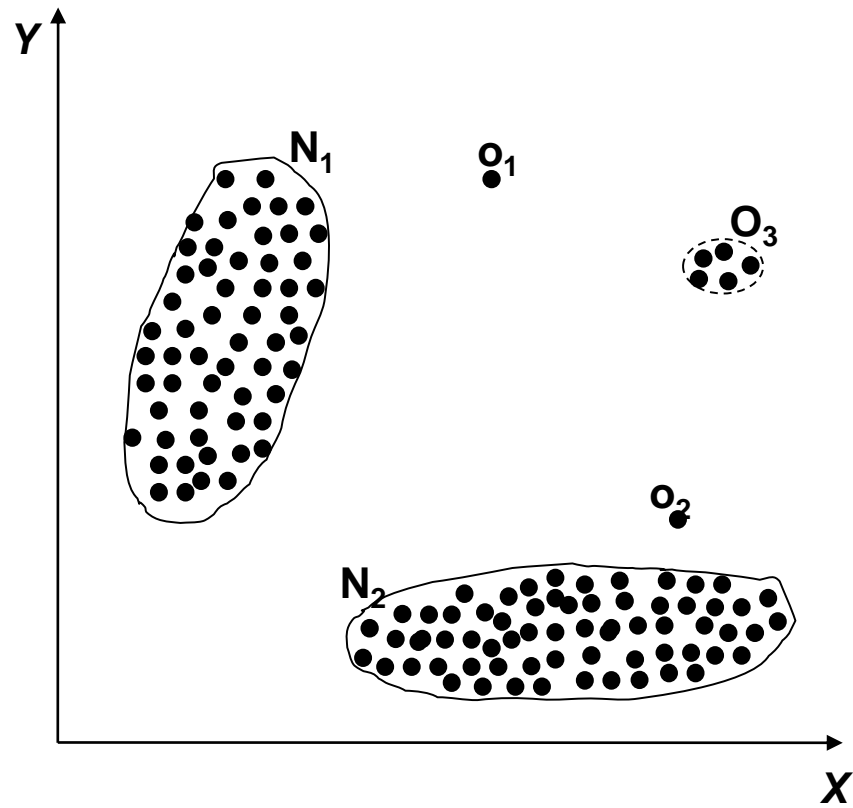
- Supervised Anomaly Detection
 - Labels available for both normal data and anomalies
 - Similar to rare class mining
- Semi-supervised Anomaly Detection
 - Labels available only for normal data
- Unsupervised Anomaly Detection
 - No labels assumed
 - Based on the assumption that anomalies are very rare compared to normal data

Type of Anomaly

- Point Anomalies
- Contextual Anomalies
- Collective Anomalies

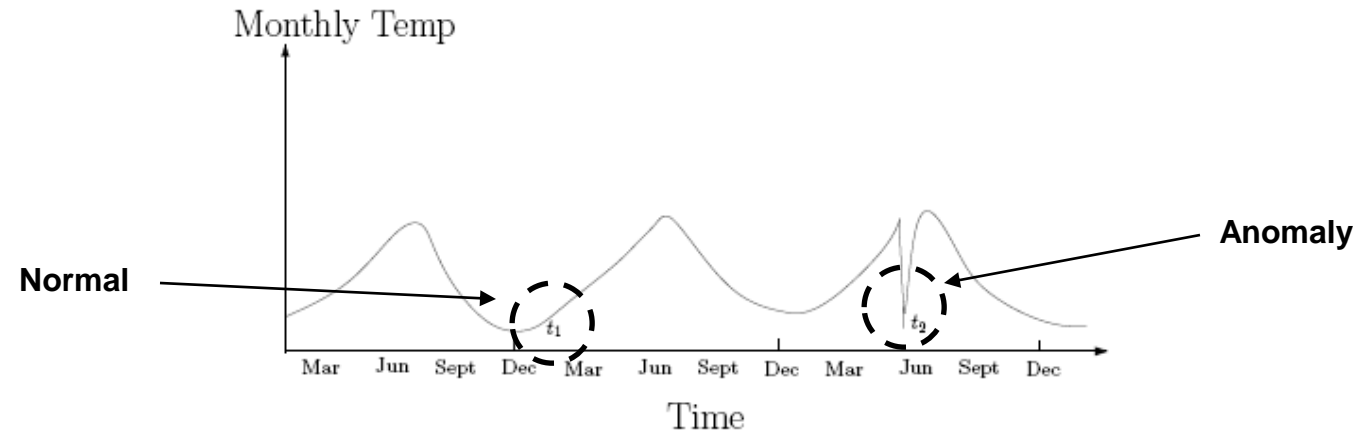
Point Anomalies

- An individual data instance is anomalous w.r.t. the data



Contextual Anomalies

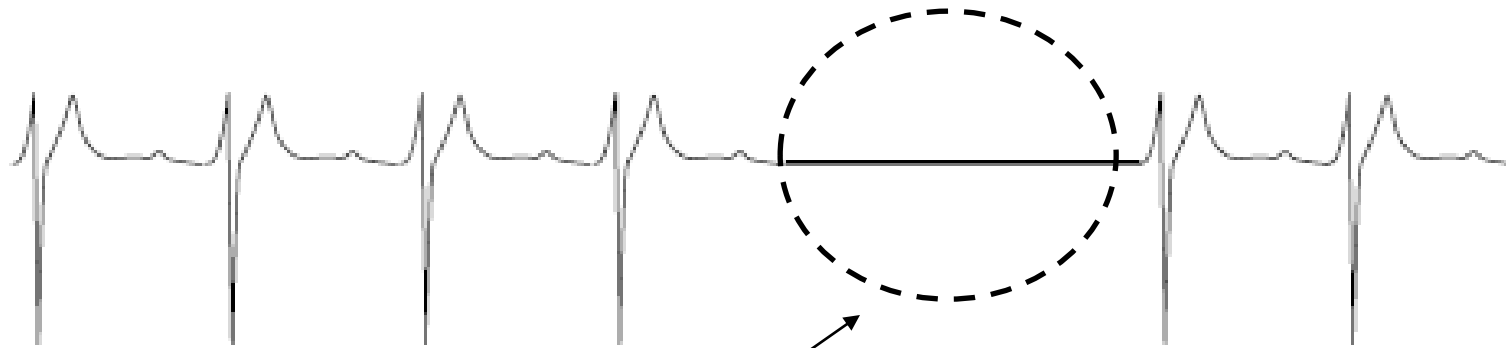
- An individual data instance is anomalous within a context
- Requires a notion of context
- Also referred to as conditional anomalies*



* Xiuyao Song, Mingxi Wu, Christopher Jermaine, Sanjay Ranka, Conditional Anomaly Detection, IEEE Transactions on Data and Knowledge Engineering, 2006.

Collective Anomalies

- A collection of related data instances is anomalous
- Requires a relationship among data instances
 - Sequential Data
 - Spatial Data
 - Graph Data
- The individual instances within a collective anomaly are not anomalous by themselves



Anomalous Subsequence

Output of Anomaly Detection

- Label
 - Each test instance is given a *normal* or *anomaly* label
 - This is especially true of classification-based approaches
- Score
 - Each test instance is assigned an anomaly score
 - Allows the output to be ranked
 - Requires an additional threshold parameter

Applications of Anomaly Detection

- Network intrusion detection
- Insurance / Credit card fraud detection
- Healthcare Informatics / Medical diagnostics
- Industrial Damage Detection
- Image Processing / Video surveillance
- Novel Topic Detection in Text Mining
- ...

Intrusion Detection

- Intrusion Detection:
 - Process of monitoring the events occurring in a computer system or network and analyzing them for intrusions
 - Intrusions are defined as attempts to bypass the security mechanisms of a computer or network
- Challenges
 - Traditional signature-based intrusion detection systems are based on signatures of known attacks and cannot detect emerging cyber threats
 - Substantial latency in deployment of newly created signatures across the computer system
- Anomaly detection can alleviate these limitations



Fraud Detection

- Fraud detection refers to detection of criminal activities occurring in commercial organizations
 - Malicious users might be the actual customers of the organization or might be posing as a customer (also known as identity theft).
- Types of fraud
 - Credit card fraud
 - Insurance claim fraud
 - Mobile / cell phone fraud
 - Insider trading
- Challenges
 - Fast and accurate real-time detection
 - Misclassification cost is very high



Healthcare Informatics

- Detect anomalous patient records
 - Indicate disease outbreaks, instrumentation errors, etc.
- Key Challenges
 - Only normal labels available
 - Misclassification cost is very high
 - Data can be complex: spatio-temporal



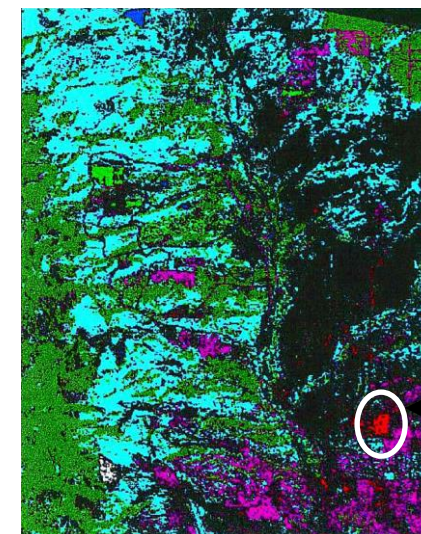
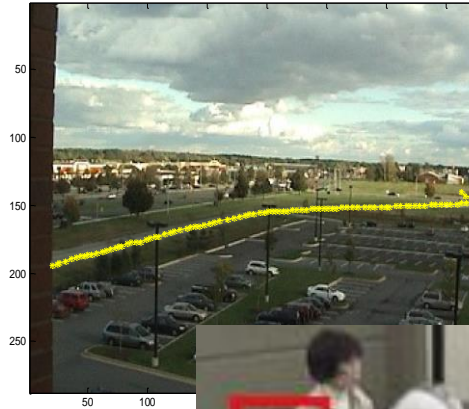
Industrial Damage Detection

- Industrial damage detection refers to detection of different faults and failures in complex industrial systems, structural damages, intrusions in electronic security systems, suspicious events in video surveillance, abnormal energy consumption, etc.
 - Example: Aircraft Safety
 - Anomalous Aircraft (Engine) / Fleet Usage
 - Anomalies in engine combustion data
 - Total aircraft health and usage management
- Key Challenges
 - Data is extremely huge, noisy and unlabelled
 - Most of applications exhibit temporal behavior
 - Detecting anomalous events typically require immediate intervention



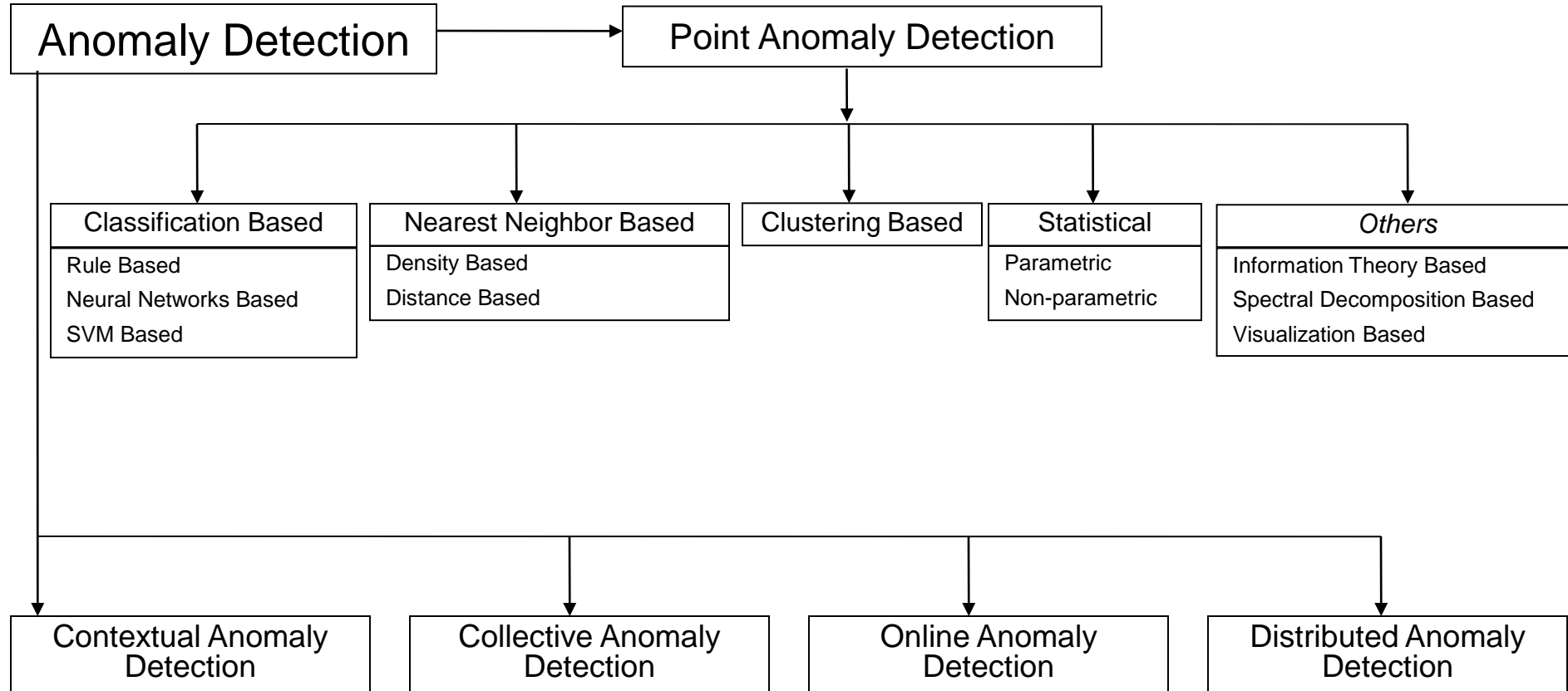
Image Processing

- Detecting outliers in a image monitored over time
- Detecting anomalous regions within an image
- Used in
 - mammography image analysis
 - video surveillance
 - satellite image analysis
- Key Challenges
 - Detecting collective anomalies
 - Data sets are very large



Anomaly

Taxonomy*



* Outlier Detection – A Survey, Varun Chandola, Arindam Banerjee, and Vipin Kumar, Technical Report TR07-17, University of Minnesota (Under Review)



THANK YOU

Munawar, PhD – moenawar@gmail.com – www.moenawar.web.id

