

Anomaly Detection

Jing Gao
SUNY Buffalo

Anomaly Detection

- **Anomalies**

- the set of objects are considerably dissimilar from the remainder of the data
- occur relatively infrequently
- when they do occur, their consequences can be quite dramatic and quite often in a negative sense



**“Mining needle in a haystack.
So much hay and so little time”**

Definition of 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

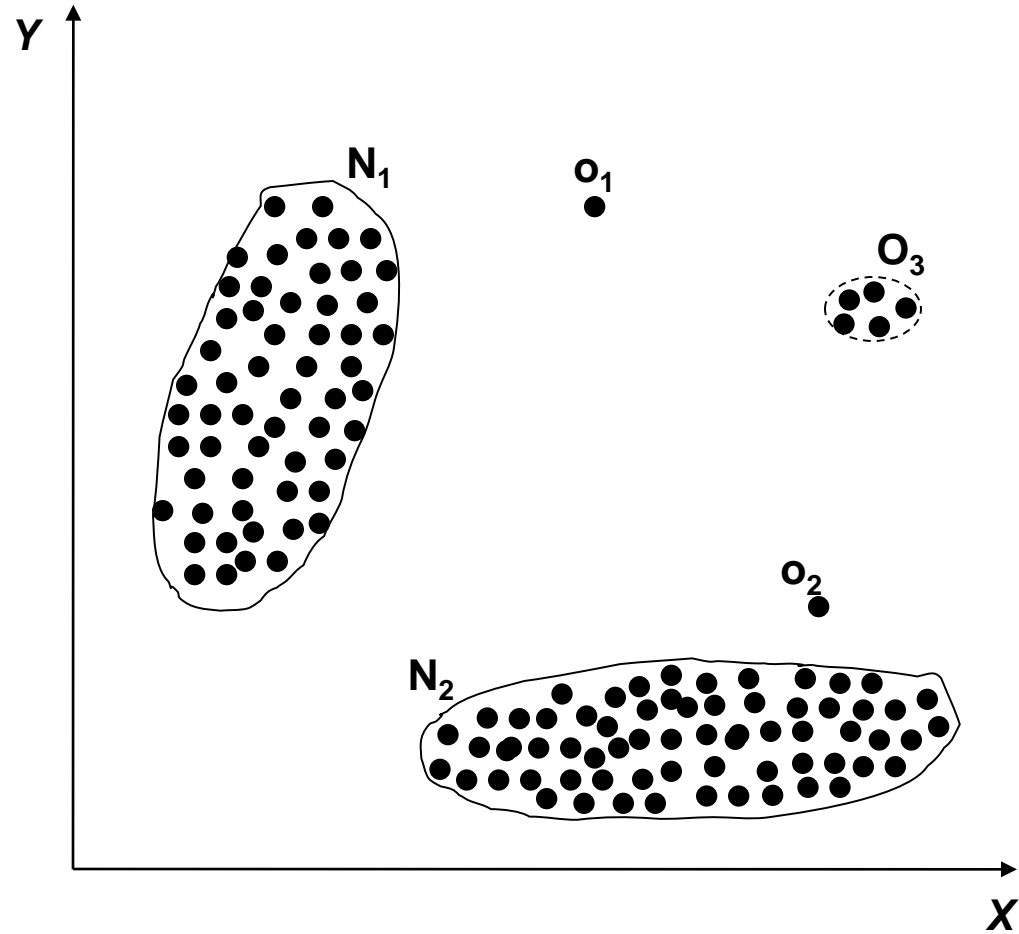
Real World Anomalies

- Credit Card Fraud
 - An abnormally high purchase made on a credit card
- Cyber Intrusions
 - Computer virus spread over Internet



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



Related problems

- Rare Class Mining
- Chance discovery
- Novelty Detection
- Exception Mining
- Noise Removal

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
- Limited availability of labeled data for training/validation
- Malicious adversaries
- Data might contain noise
- Normal behaviour 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

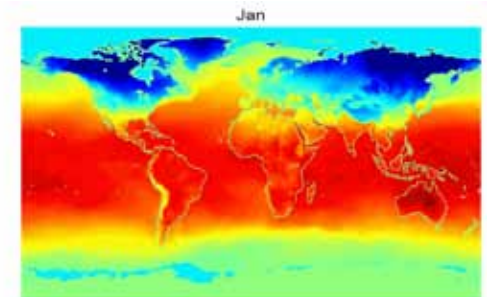
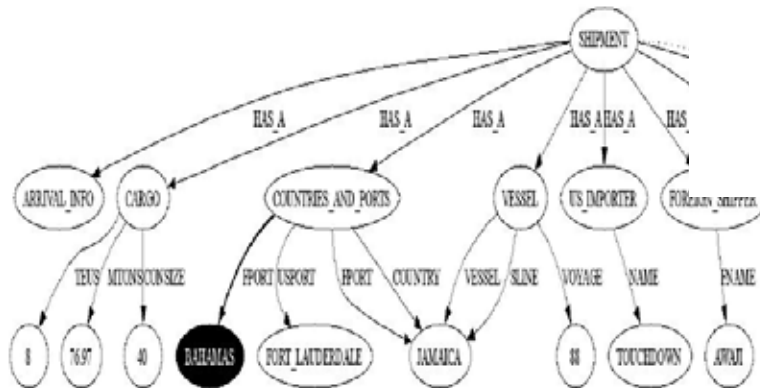
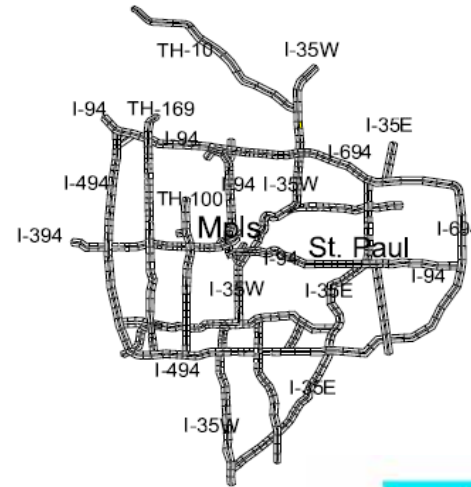
- Most common form of data handled by anomaly detection techniques is *Record Data*
 - Univariate
 - Multivariate

<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 – Complex Data Types

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

```
GGTTCCGCCTTCAGCCCCGCGCC
CGCAGGGCCCGCCCCGCGCCGTC
GAGAAGGGCCCGCCTGGCGGGCG
GGGGGAGGCGGGGCCGCCCGAGC
CCAACCGAGTCCGACCAGGTGCC
CCCTCTGCTCGGCCTAGACCTGA
GCTCATTAGGCGGCAGCGGACAG
GCCAAGTAGAACACGCGAAGCGC
TGGGCTGCTGCTGCGACCAGGG
```



Data Labels

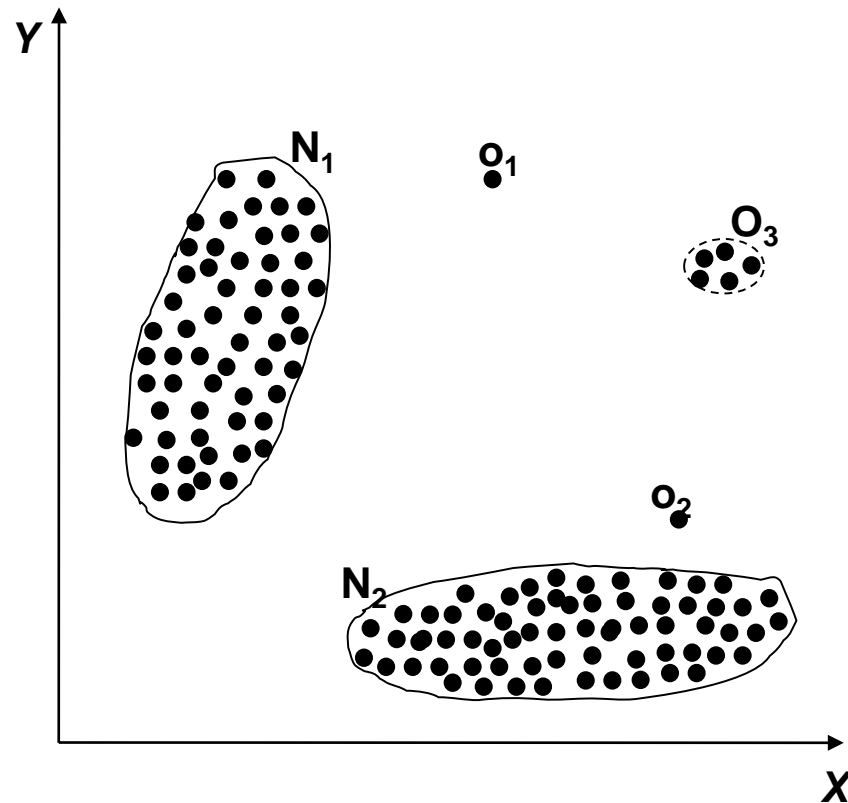
- **Supervised Anomaly Detection**
 - Labels available for both normal data and anomalies
 - Similar to skewed (imbalanced) classification
- **Semi-supervised Anomaly Detection**
 - Limited amount of labeled data
 - Combine supervised and unsupervised techniques
- **Unsupervised Anomaly Detection**
 - No labels assumed
 - Based on the assumption that anomalies are very rare compared to normal data

Type of Anomalies

- Point Anomalies
- Contextual Anomalies
- Collective Anomalies

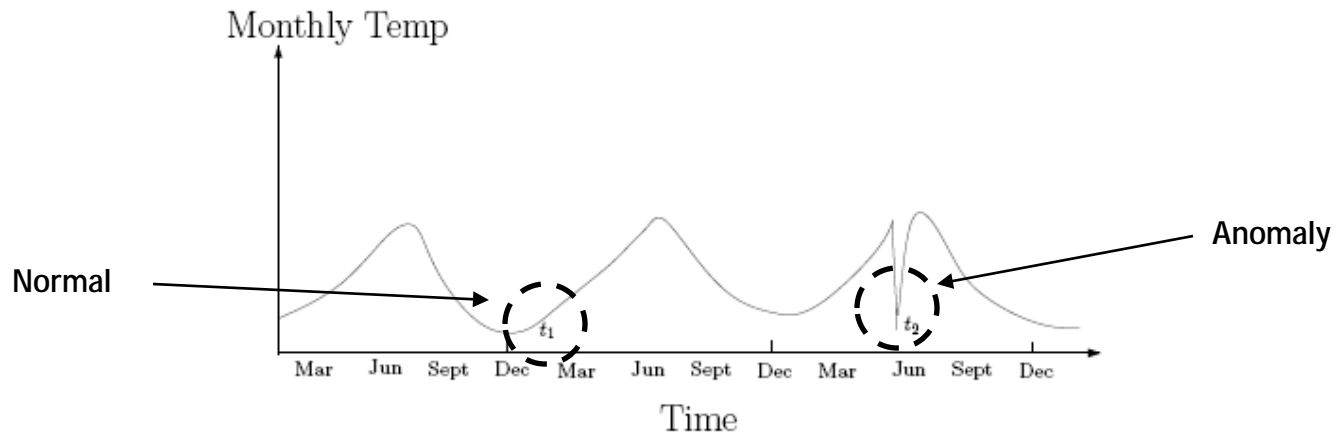
Point Anomalies

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



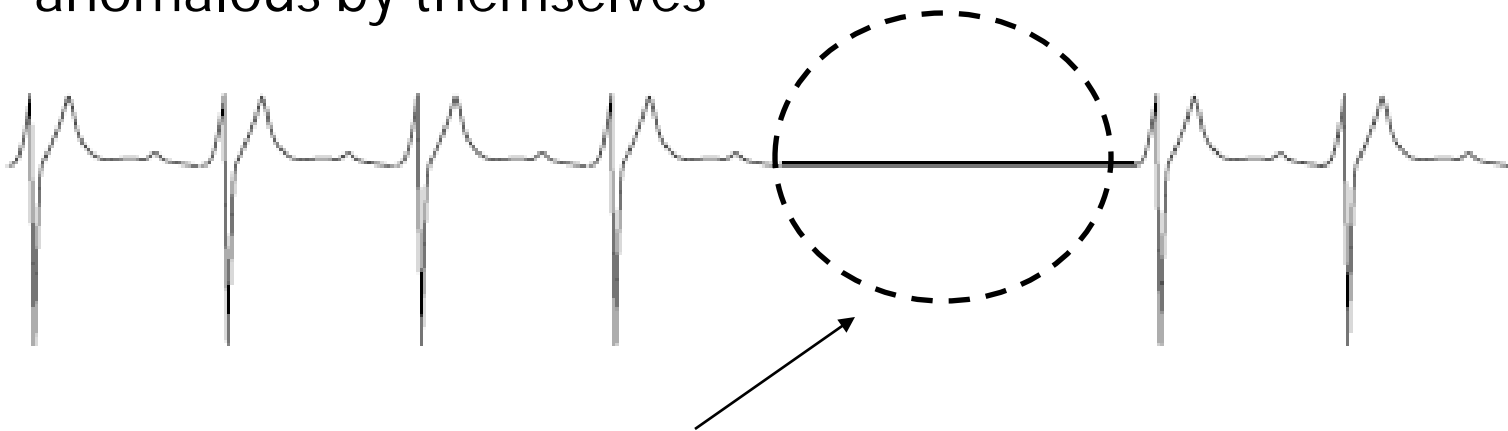
Contextual Anomalies

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



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

Metrics for Performance Evaluation

- Confusion Matrix

	PREDICTED CLASS		
		+	-
ACTUAL CLASS	+	a	b
	-	c	d

a: TP (true positive)

c: FP (false positive)

b: FN (false negative)

d: TN (true negative)

Metrics for Performance Evaluation

	PREDICTED CLASS		
	+	-	
ACTUAL CLASS	+	a (TP)	b (FN)
	-	c (FP)	d (TN)

- Measure used in classification:

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

Limitation of Accuracy

- Anomaly detection
 - Number of negative examples = 9990
 - Number of positive examples = 10
- If model predicts everything to be class 0, accuracy is $9990/10000 = 99.9\%$
 - Accuracy is misleading because model does not detect any positive examples

Cost Matrix

	PREDICTED CLASS		
ACTUAL CLASS	$C(i j)$	+	-
	+	$C(+ +)$	$C(- +)$
	-	$C(+ -)$	$C(- -)$

$C(i|j)$: Cost of misclassifying class j example as class i

Computing Cost of Classification

Cost Matrix	PREDICTED CLASS		
	C(i j)	+	-
ACTUAL CLASS	+	-1	100
	-	1	0

Model M_1	PREDICTED CLASS		
		+	-
ACTUAL CLASS	+	150	40
	-	60	250

Accuracy = 80%

Cost = 3910

Model M_2	PREDICTED CLASS		
		+	-
ACTUAL CLASS	+	250	45
	-	5	200

Accuracy = 90%

Cost = 4255

Cost-Sensitive Measures

$$\text{Precision (p)} = \frac{a}{a + c}$$

$$\text{Recall (r)} = \frac{a}{a + b}$$

$$\text{F - measure (F)} = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$$

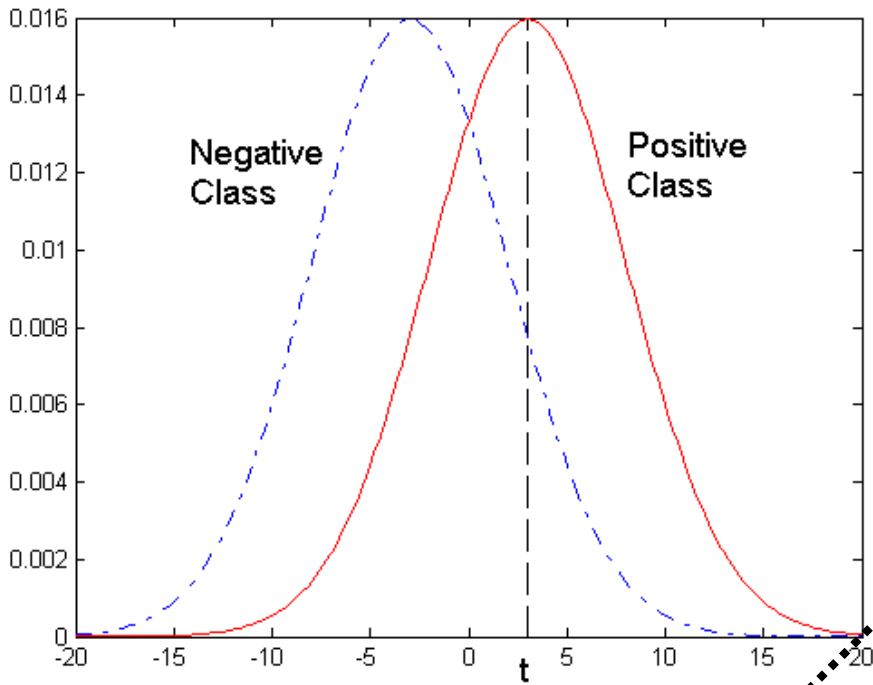
$$\text{Weighted Accuracy} = \frac{w_1 a + w_4 d}{w_1 a + w_2 b + w_3 c + w_4 d}$$

ROC (Receiver Operating Characteristic)

- ROC curve plots TPR (on the y-axis) against FPR (on the x-axis)
- Performance of each classifier represented as a point on the ROC curve
 - changing the threshold of algorithm, sample distribution or cost matrix changes the location of the point

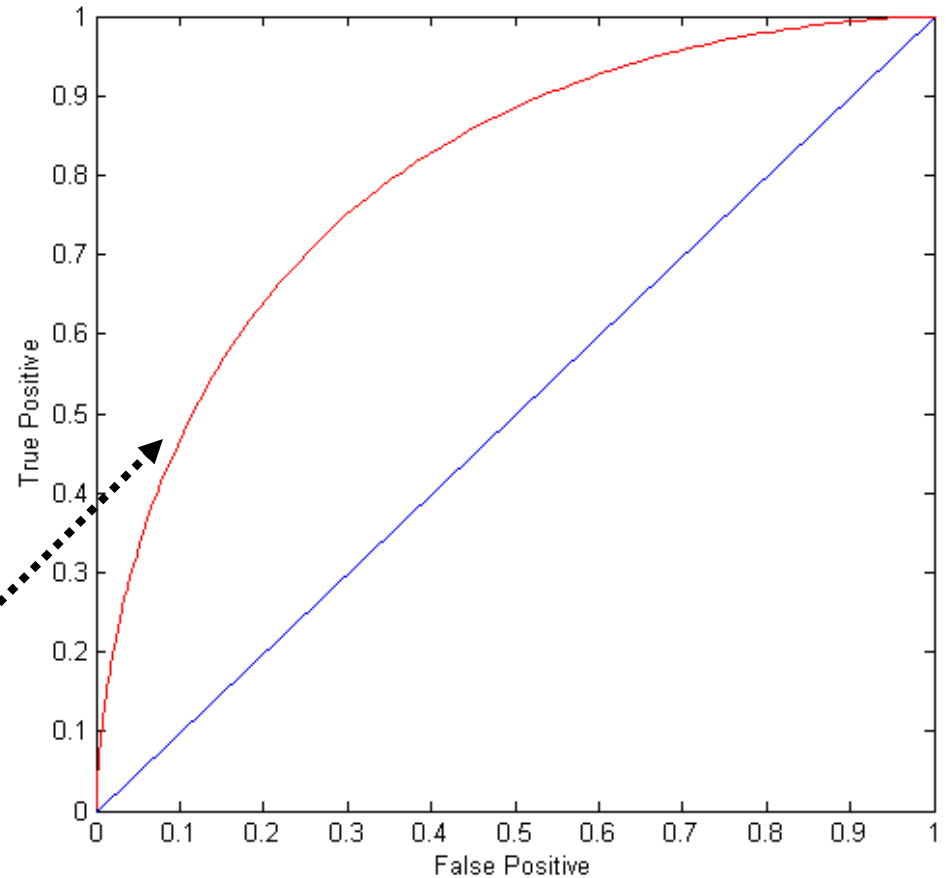
ROC Curve

- 1-dimensional data set containing 2 classes (positive and negative)
- any points located at $x > t$ is classified as positive



At threshold t :

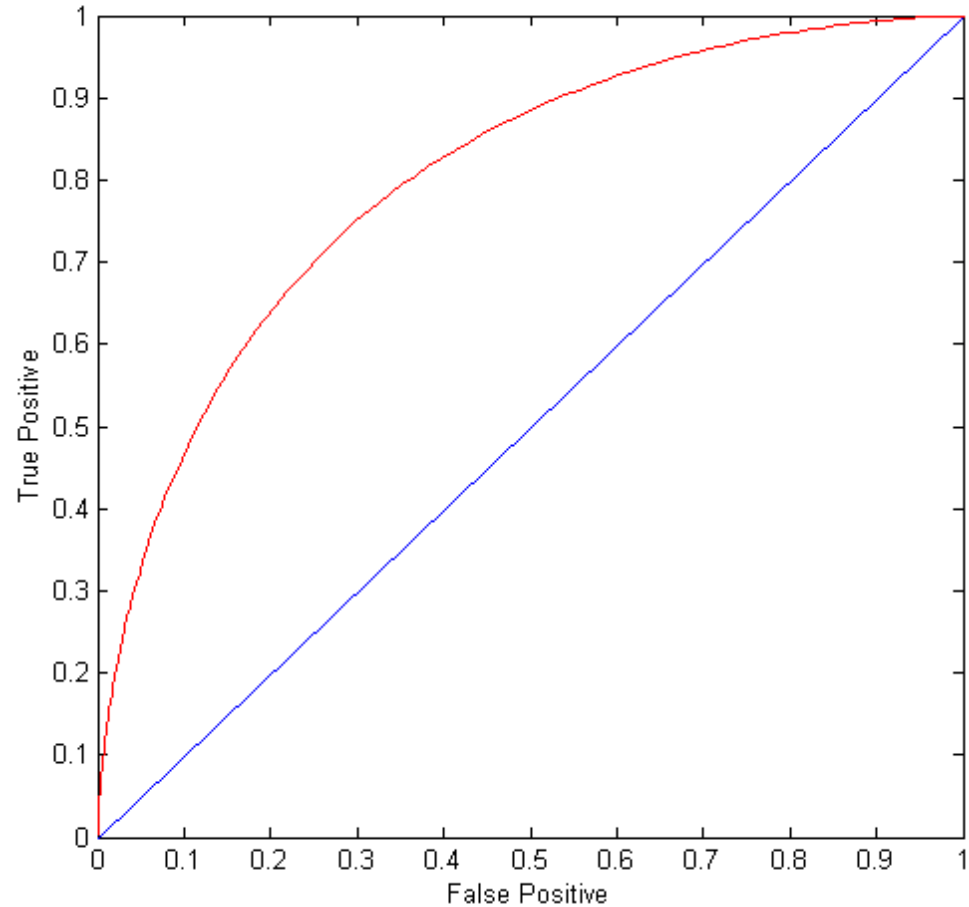
TP=0.5, FN=0.5, FP=0.12, FN=0.88



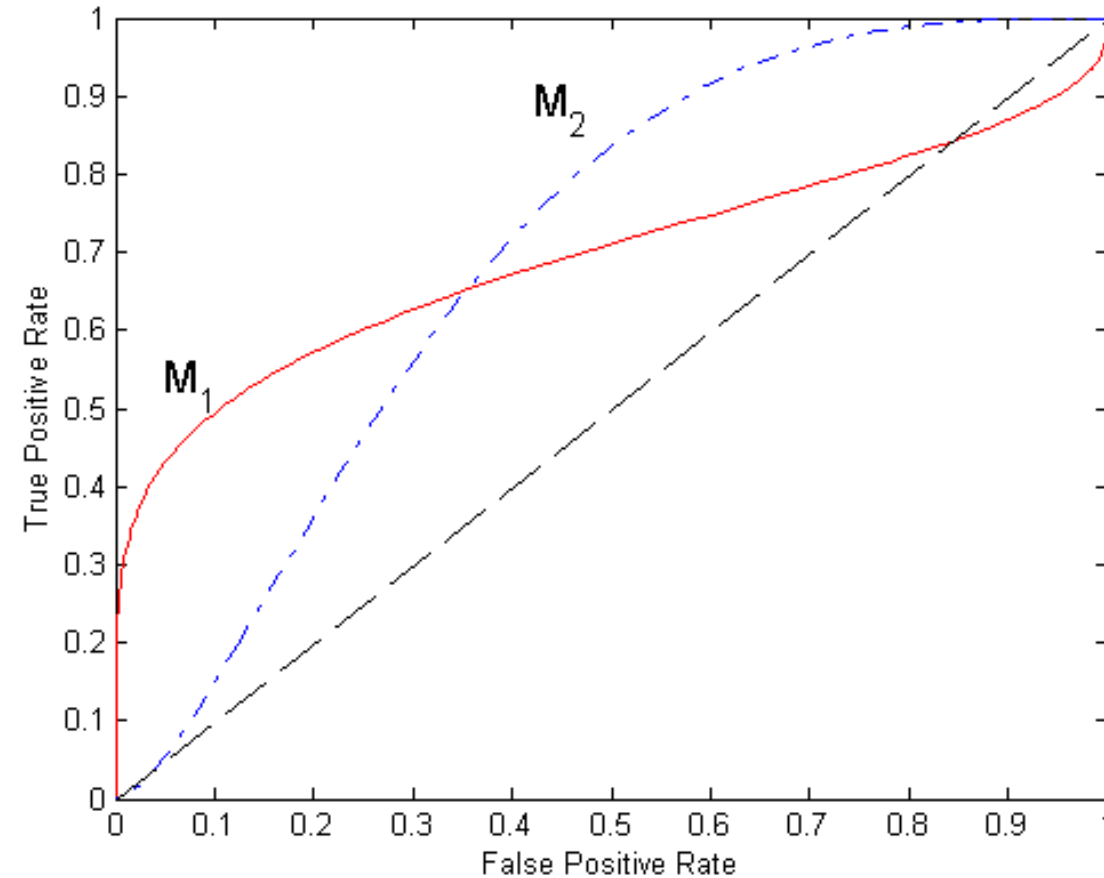
ROC Curve

(TPR, FPR):

- (0,0): declare everything to be negative class
- (1,1): declare everything to be positive class
- (1,0): ideal
- Diagonal line:
 - Random guessing
 - Below diagonal line:
 - prediction is opposite of true class



Using ROC for Model Comparison



- Comparing two models
 - M₁ is better for small FPR
 - M₂ is better for large FPR
- Area Under the ROC curve
 - Ideal:
 - § Area = 1
 - Random guess:
 - § Area = 0.5

How to Construct an ROC curve

Instance	Score	Label
1	0.95	+
2	0.93	+
3	0.87	-
4	0.85	-
5	0.85	-
6	0.85	+
7	0.76	-
8	0.53	+

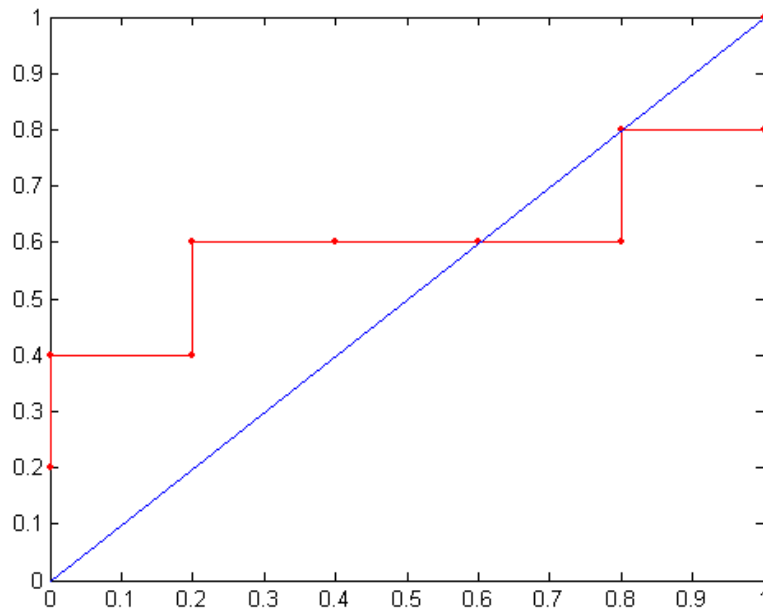
	PREDICTED CLASS		
		+	-
ACTUAL CLASS	+	a (TP)	b (FN)
	-	c (FP)	d (TN)

- Calculate the outlier scores of the given instances
- Sort the instances according to the scores in decreasing order
- Apply threshold at each unique value of the score
- Count the number of TP, FP, TN, FN at each threshold
- TP rate, $TPR = TP/(TP+FN)$
- FP rate, $FPR = FP/(FP + TN)$

How to construct an ROC curve

Class	+	-	+	-	-	-	+	-	+	+	
Threshold \geq	0.25	0.43	0.53	0.76	0.85	0.85	0.85	0.87	0.93	0.95	1.00
TP	5	4	4	3	3	3	3	2	2	1	0
FP	5	5	4	4	3	2	1	1	0	0	0
TN	0	0	1	1	2	3	4	4	5	5	5
FN	0	1	1	2	2	2	2	3	3	4	5
TPR	1	0.8	0.8	0.6	0.6	0.6	0.6	0.4	0.4	0.2	0
FPR	1	1	0.8	0.8	0.6	0.4	0.2	0.2	0	0	0

ROC Curve:



Applications of Anomaly Detection

- Network intrusion detection
- Insurance / Credit card fraud detection
- Healthcare Informatics / Medical diagnostics
- Image Processing / Video surveillance
- ...

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



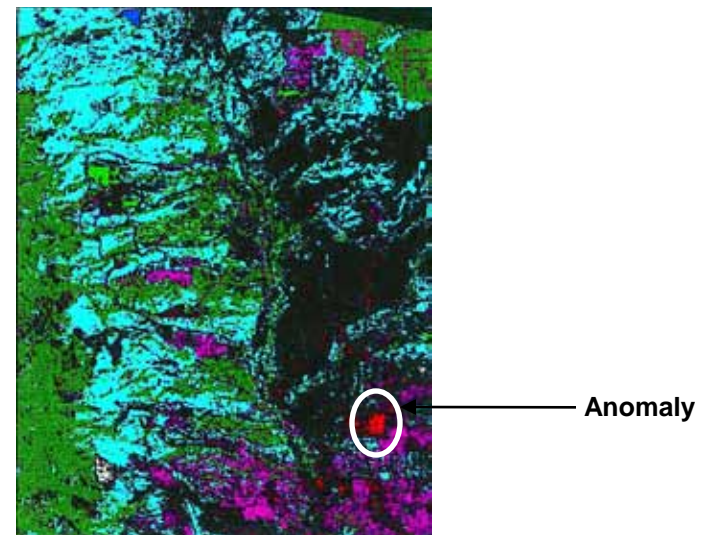
Healthcare Informatics

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



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



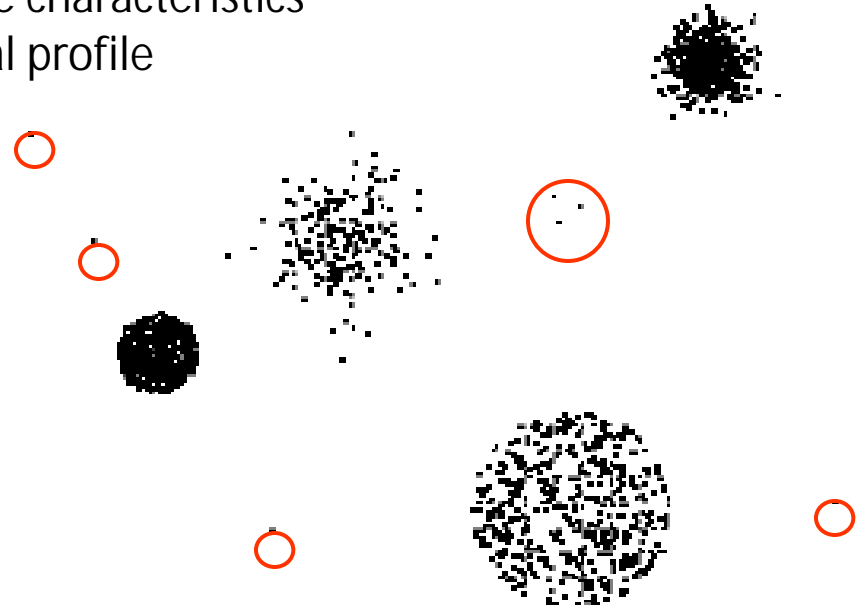
Anomaly Detection Schemes

- **General Steps**

- Build a profile of the “normal” behavior
 - Profile can be patterns or summary statistics for the overall population
- Use the “normal” profile to detect anomalies
 - Anomalies are observations whose characteristics differ significantly from the normal profile

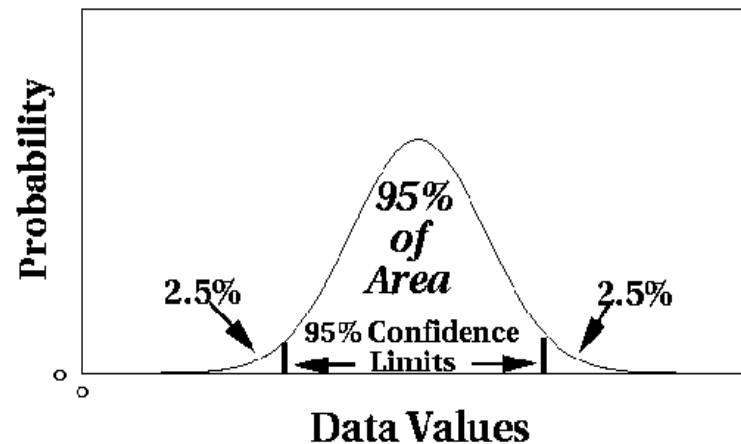
- **Methods**

- Statistical-based
- Distance-based
- Model-based



Statistical Approaches

- Assume a parametric model describing the distribution of the data (e.g., normal distribution)
- Apply a statistical test that depends on
 - Data distribution
 - Parameter of distribution (e.g., mean, variance)
 - Number of expected outliers (confidence limit)



Grubbs' Test

- Detect outliers in univariate data
- Assume data comes from normal distribution
- Detects one outlier at a time, remove the outlier, and repeat
 - H_0 : There is no outlier in data
 - H_A : There is at least one outlier

- Grubbs' test statistic:
$$G = \frac{\max |X - \bar{X}|}{s}$$

- Reject H_0 if:
$$G > \frac{(N - 1)}{\sqrt{N}} \sqrt{\frac{t^2_{(a/N, N-2)}}{N - 2 + t^2_{(a/N, N-2)}}}$$

Statistical-based – Likelihood Approach

- Assume the data set D contains samples from a mixture of two probability distributions:
 - M (majority distribution)
 - A (anomalous distribution)
- General Approach:
 - Initially, assume all the data points belong to M
 - Let $L_t(D)$ be the log likelihood of D at time t
 - For each point x_t that belongs to M , move it to A
 - Let $L_{t+1}(D)$ be the new log likelihood.
 - Compute the difference, $D = L_t(D) - L_{t+1}(D)$
 - If $D > c$ (some threshold), then x_t is declared as an anomaly and moved permanently from M to A

Statistical-based – Likelihood Approach

- Data distribution, $D = (1 - I) M + I A$
- M is a probability distribution estimated from data
 - Can be based on any modeling method, e.g., mixture model
- A can be assumed to be uniform distribution
- Likelihood at time t :

$$L_t(D) = \prod_{i=1}^N P_D(x_i) = (1 - I)^{|M_t|} \prod_{x_i \in M_t} P_{M_t}(x_i) / |A_t| \prod_{x_i \in A_t} P_{A_t}(x_i)$$

Limitations of Statistical Approaches

- Most of the tests are for a single attribute
- In many cases, data distribution may not be known
- For high dimensional data, it may be difficult to estimate the true distribution

Distance-based Approaches

- Data is represented as a vector of features
- Three major approaches
 - Nearest-neighbor based
 - Density based
 - Clustering based

Nearest-Neighbor Based Approach

- **Approach:**
 - Compute the distance between every pair of data points
 - There are various ways to define outliers:
 - Data points for which there are fewer than p neighboring points within a distance D
 - The top n data points whose distance to the k -th nearest neighbor is greatest
 - The top n data points whose average distance to the k nearest neighbors is greatest

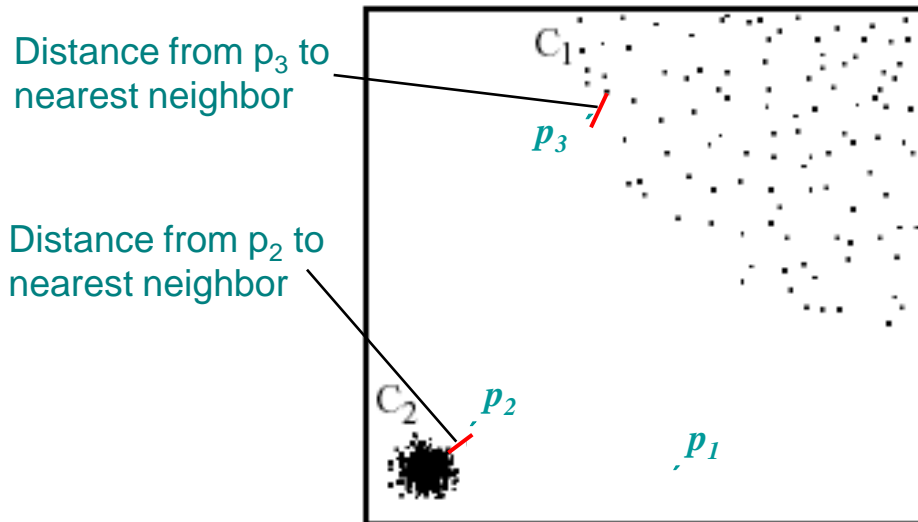
Distance-Based Outlier Detection

- For each object o , examine the # of other objects in the r -neighborhood of o , where r is a user-specified **distance threshold**
- An object o is an outlier if most (taking π as a **fraction threshold**) of the objects in D are far away from o , i.e., not in the r -neighborhood of o
- An object o is a $DB(r, \pi)$ outlier if
$$\frac{\|\{o' | dist(o, o') \leq r\}\|}{\|D\|} \leq \pi$$
- Equivalently, one can check the distance between o and its k -th nearest neighbor o_k , where $k = \lceil \pi \|D\| \rceil$. o is an outlier if $dist(o, o_k) > r$

Density-based Approach

- Local Outlier Factor (LOF) approach

– Example:



In the *NN* approach, p_2 is not considered as outlier, while the *LOF* approach find both p_1 and p_2 as outliers

NN approach may consider p_3 as outlier, but *LOF* approach does not

Density-based: LOF approach

- For each point, compute the density of its local neighborhood
- Compute local outlier factor (LOF) of a sample p as the average of the ratios of the density of sample p and the density of its nearest neighbors
- Outliers are points with largest LOF value

Local Outlier Factor: LOF

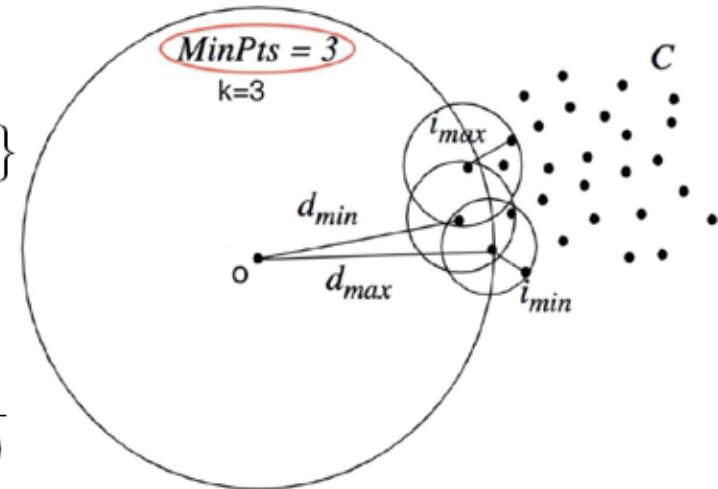
- Reachability distance from o' to o :

$$reachdist_k(o \leftarrow o') = \max\{dist_k(o), dist(o, o')\}$$

- where k is a user-specified parameter

- Local reachability density of o :

$$lrd_k(o) = \frac{\|N_k(o)\|}{\sum_{o' \in N_k(o)} reachdist_k(o' \leftarrow o)}$$



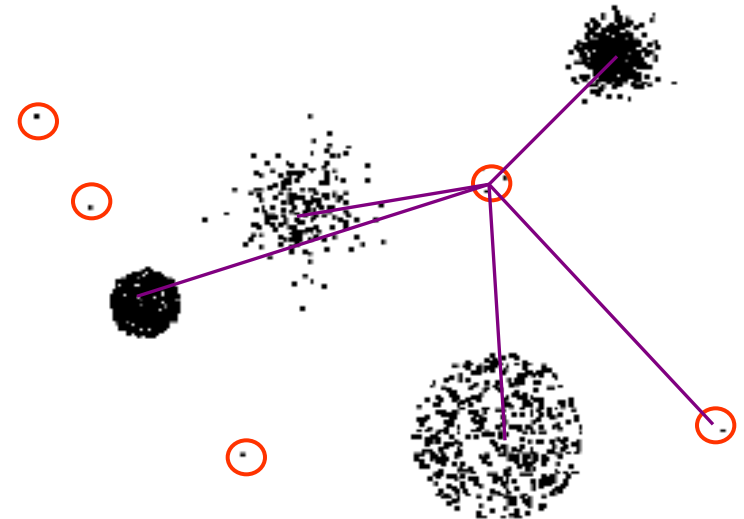
- n LOF (Local outlier factor) of an object o is the average of the ratio of local reachability of o and those of o 's k -nearest neighbors

$$LOF_k(o) = \frac{\sum_{o' \in N_k(o)} \frac{lrd_k(o')}{lrd_k(o)}}{\|N_k(o)\|} = \sum_{o' \in N_k(o)} lrd_k(o') \cdot \sum_{o' \in N_k(o)} reachdist_k(o' \leftarrow o)$$

- n The higher the local reachability distance of o , and the higher the local reachability density of the k NN of o , the higher LOF
- n This captures a local outlier whose local density is relatively low comparing to the local densities of its k NN

Clustering-Based

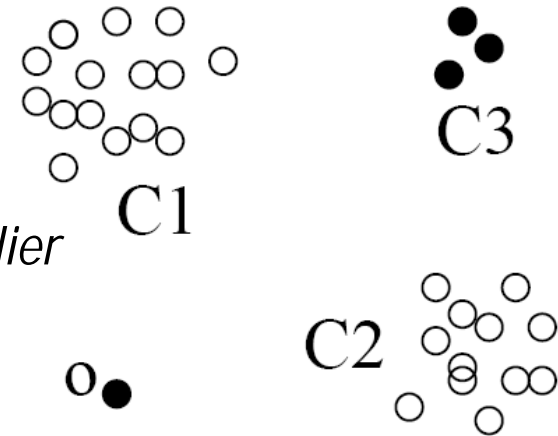
- **Basic idea:**
 - Cluster the data into groups of different density
 - Choose points in small cluster as candidate outliers
 - Compute the distance between candidate points and non-candidate clusters.
 - If candidate points are far from all other non-candidate points, they are outliers



Detecting Outliers in Small Clusters

- *FindCBLOF*: Detect outliers in small clusters

- Find clusters, and sort them in decreasing size
- To each data point, assign a *cluster-based local outlier factor* (CBLOF):
 - If obj p belongs to a large cluster, $CBLOF = \text{cluster_size} \times \text{similarity between } p \text{ and cluster}$
 - If p belongs to a small one, $CBLOF = \text{cluster size} \times \text{similarity betw. } p \text{ and the closest large cluster}$



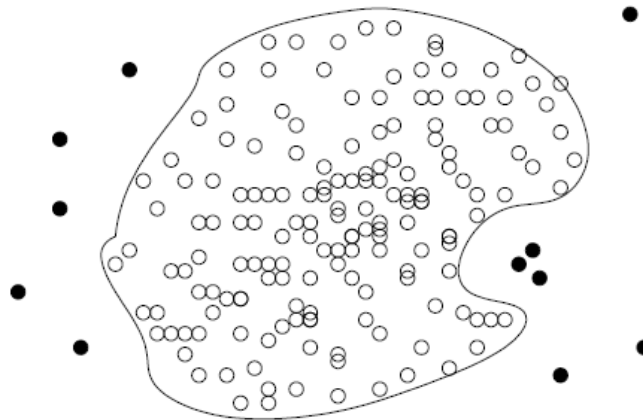
- n Ex. In the figure, o is outlier since its closest large cluster is C_1 , but the similarity between o and C_1 is small. For any point in C_3 , its closest large cluster is C_2 but its similarity from C_2 is low, plus $|C_3| = 3$ is small

Classification-Based Methods

- Idea: Train a classification model that can distinguish “normal” data from outliers
- Consider a training set that contains samples labeled as “normal” and others labeled as “outlier”
 - But, the training set is typically heavily biased: # of “normal” samples likely far exceeds # of outlier samples
- Handle the imbalanced distribution
 - Oversampling positives and/or undersampling negatives
 - Alter decision threshold
 - Cost-sensitive learning

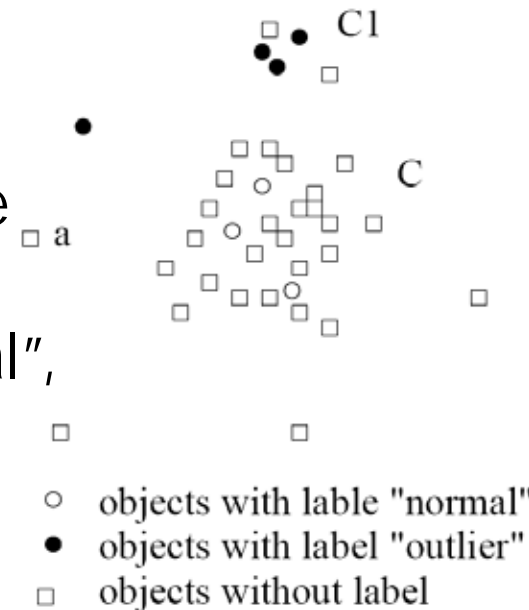
One-Class Model

- n One-class model: A classifier is built to describe only the normal class
 - n Learn the decision boundary of the normal class using classification methods such as SVM
 - n Any samples that do not belong to the normal class (not within the decision boundary) are declared as outliers
 - n Adv: can detect new outliers that may not appear close to any outlier objects in the training set



Semi-Supervised Learning

- Semi-supervised learning: Combining classification-based and clustering-based methods
- Method
 - Using a clustering-based approach, find a large cluster, C , and a small cluster, C_1
 - Since some objects in C carry the label "normal", treat all objects in C as normal
 - Use the one-class model of this cluster to identify normal objects in outlier detection
 - Since some objects in cluster C_1 carry the label "outlier", declare all objects in C_1 as outliers
 - Any object that does not fall into the model for C (such as a) is considered an outlier as well



Take-away Message

- Definition of outlier detection
- Applications of outlier detection
- Evaluation of outlier detection techniques
- Unsupervised approaches (statistical, distance, density-based)
- Supervised and semi-supervised approaches