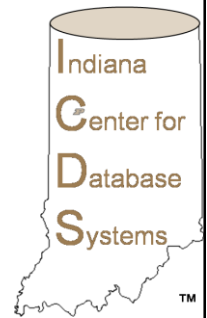


CS57300: Data Mining

Anomaly Detection

Prof. Chris Clifton

14 April 2022



Some materials from *Introduction to Data Mining* by Tan, Steinbach and Kumar

Task

- **Anomalies/outliers:** data points that are considerably “different” from the remainder of the data
- Variants:
 - Find all points with anomaly scores $>$ threshold
 - Find point with largest anomaly score
 - Given a database D with mostly normal points, compute the anomaly score of a point x with respect to D

Examples

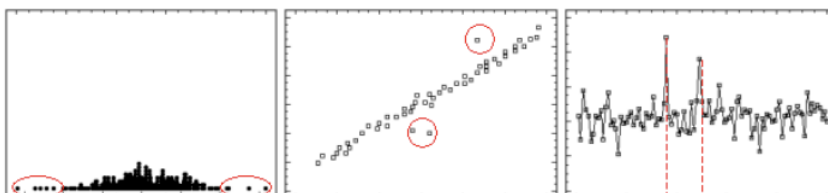
- Fraud detection
- Intrusion detection
- Ecosystem disturbances
- System monitoring
- Biosurveillance/public health
- Data preprocessing

Types of anomalies

- Data from different classes
 - “An outlier is an observation that differs so much from other observations as to arouse suspicion that it was generated by a different mechanism”
- Natural variation
 - Extreme or unlikely variations are often interesting
- Data measurement and collection errors
 - Preprocess to remove

Defining an outlier

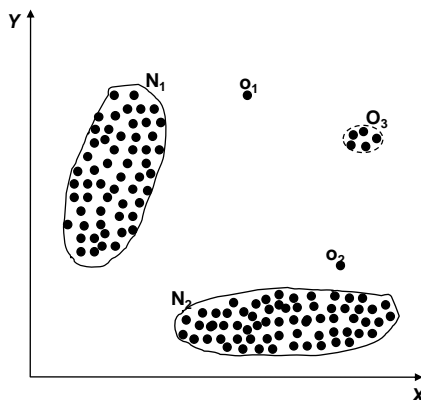
- Notion of outlier is highly subjective and domain dependent
- However, most definitions can be viewed as defining a distribution for “normal” data and then looking for deviations from that distribution



Source: Osmar Zaiane, UAlberta, PKDD

Point anomalies

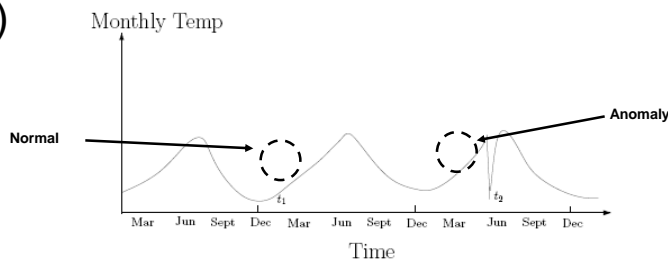
- An **individual** data instance is anomalous with respect to the **data**



Source: Lazarevic et al, ECML/PKDD'08 Tutorial

Contextual anomalies

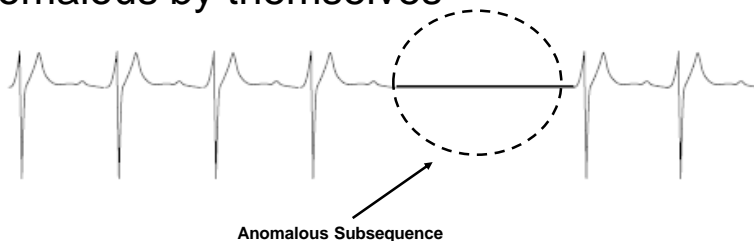
- An **individual** data instance is anomalous within a **context**
- Requires a notion of context
- Also referred to as conditional anomalies (Song et. al, TDKE '06)



Source: Lazarevic et al, ECML/PKDD'08 Tutorial

Collective anomalies

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



Source: Lazarevic et al, ECML/PKDD'08 Tutorial

Anomaly detection

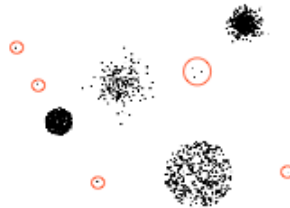
- Challenges
 - How many attributes are used to define an outlier?
 - How many outliers are there in the data?
 - Class labels are costly (evaluation can be challenging)
 - Skewed class distribution (finding needles in haystack)
- Working assumption:
 - There are considerably more “normal” observations than “abnormal” observations in the data

Approaches

- Supervised
 - Labels available for both normal data and anomalies
 - Similar to classification with imbalanced classes
- Semi-supervised
 - Labels available only for normal data
- Unsupervised
 - No labels assumed
 - Based on the assumption that anomalies are very rare compared to normal data

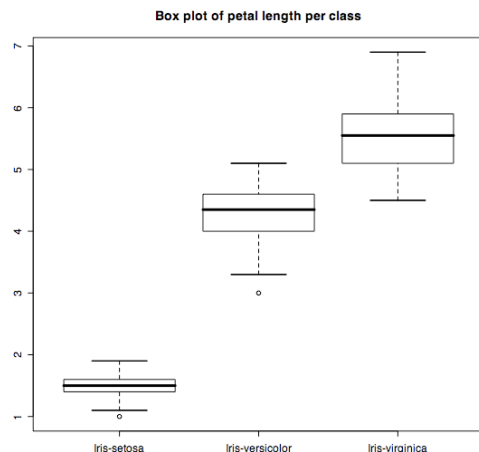
Unsupervised (point) anomaly detection

- General method
 - Build a profile of “normal” behavior based on patterns or summary statistics for the overall population
 - Use deviations from “normal” to detect anomalies
- Types of methods
 - Visual and statistical-based
 - Distance-based
 - Model-based



Visual methods

- Box plot (1D)
- Scatter plot (2D)
- Limitations
 - Time consuming
 - Subjective



Distance-based approaches

- Three major types of methods
 - Nearest-neighbor
 - Density-based
 - Clustering approach

Nearest-neighbor

- Compute distance between every pair of points
- How to define outliers?
 - Points for which there are fewer than p neighboring points within distance d
 - Top p points whose distance to k^{th} nearest neighbor is greatest
 - Top p points whose average distance to their k nearest neighbors is greatest

Example

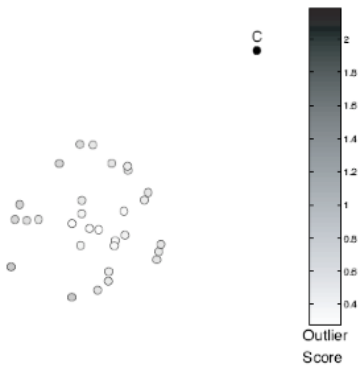


Figure 10.4. Outlier score based on the distance to the **fifth** nearest neighbor.

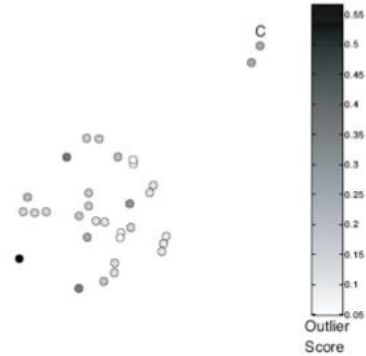


Figure 10.5. Outlier score based on the distance to the **first** nearest neighbor. Nearby outliers have low outlier scores.

Example

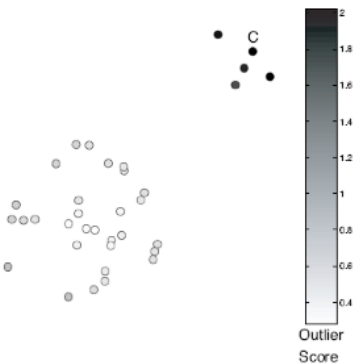


Figure 10.6. Outlier score based on distance to the **fifth** nearest neighbor. A small cluster becomes an outlier.

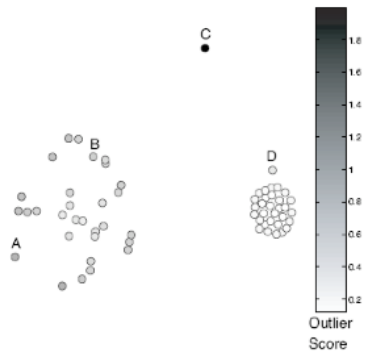


Figure 10.7. Outlier score based on the distance to the **fifth** nearest neighbor. Clusters of differing density.

Density-based

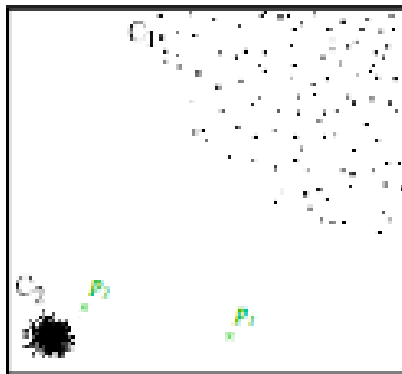
- For each point, compute the density of its local neighborhood of k neighbors

$$\text{density}(x, k) = \frac{|N(x, k)|}{\sum_{y \in N(x, k)} \text{distance}(x, y)}$$

- Local outlier factor (LOF) is the ratio of a point's density to the average density of its nearest neighbors

$$\text{LOF}(x) = \frac{\frac{1}{k} \sum_{y \in N(x, k)} \text{density}(y, k)}{\text{density}(x, k)}$$

- Outliers are points with largest LOF value



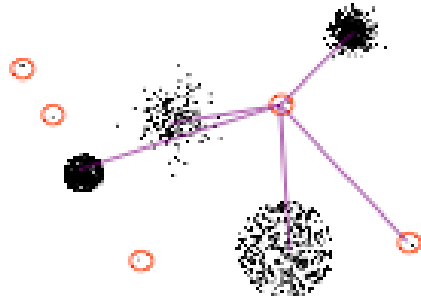
In the NN approach, p_2 is not discovered as an outlier, while the LOF approach considers both p_1 and p_2 to be outliers

High dimensions

- In high-dimensional space, data is sparse and notion of proximity becomes meaningless
 - Every point is almost equally good outlier from the perspective of proximity-based definitions
- Lower-dimensional projections can be used for outlier detection
 - A point is an outlier if in some lower dimensional projection, it is present in a local region of abnormally low density

Clustering-based

- 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



Example

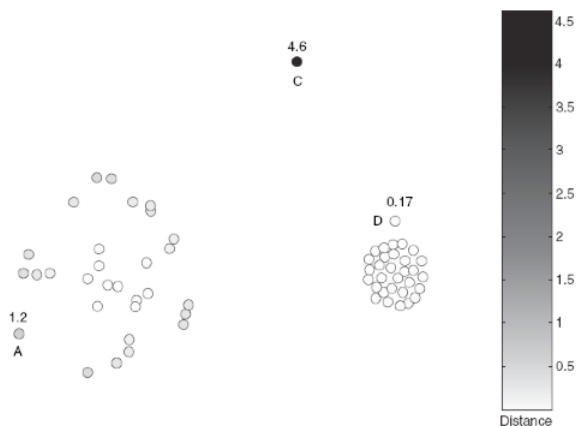


Figure 10.9. Distance of points from closest centroid.

This doesn't take into account the average distance of points to their cluster centroid (which can vary by cluster density)... so use relative distance (ratio of distance to median distance)

Example

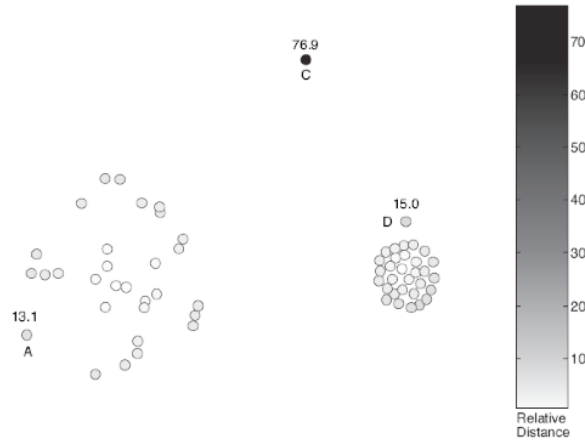


Figure 10.10. Relative distance of points from closest centroid.

Supervised Anomaly Detection

- *Anomalies* not known in advance
 - Otherwise they wouldn't be anomalies
- But what if we assume we know normal?
 - Training data is from non-anomalies
- Train classifier to recognize normal
- One-Class Classification (*Moya & Hush '96*)

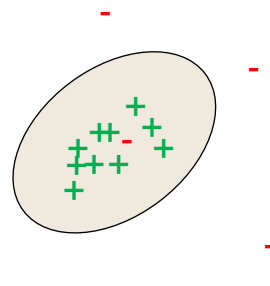
One-class Classification

- Problem: “Easy” to classify training data
 - Given instance, classify as normal
 - Works for all the training data
- Ideas:
 - “Fake” training data
 - Unpruned classifier
 - Narrowly tailor to recognize positive instances

23

Artificial Training Data

- Start with a large volume of “normal” data
- Randomly generate a comparatively small volume of presumed anomalies
 - From uniform distribution across space
- Classifier accepts some false positives

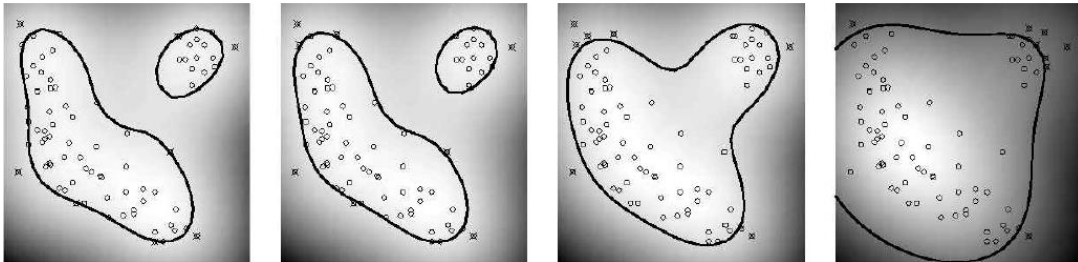


24

One-Class Support Vector Machine

Schölkopf, Platt, Shawe-Taylor, Smola and Williamson 2001

- Idea: “Max-Margin” between points and “everything else”
 - Accomplished through clever choice of kernel functions
 - Max margin between points and origin gives boundary separating points from “everything else”

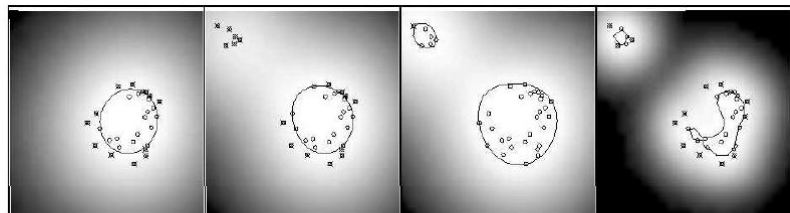


26

One-Class Support Vector Machine

Schölkopf, Platt, Shawe-Taylor, Smola and Williamson 2001

- Parameters
 - ν bounds allowed outliers
 - γ governs classifier complexity (kernel function space)



ν , width c	0.5, 0.5	0.5, 0.5	0.1, 0.5	0.5, 0.1
frac. SVs/OLs	0.54, 0.43	0.59, 0.47	0.24, 0.03	0.65, 0.38
margin $\rho/\ w\ $	0.84	0.70	0.62	0.48

27

Anomaly Detection: Statistical

Some materials from *Introduction to Data Mining* by Tan, Steinbach and Kumar

Statistical approaches

- Use a parametric model to describe the data (e.g., Normal distribution)
- Apply a statistical test that evaluates how likely a point is under the data distribution
- Need to specify a confidence limit (e.g., 3σ away from mean is an outlier)

Multivariate data

1. Statistical approach

- Model data with a multivariate Gaussian
- Calculate the likelihood of a point with respect to the estimated distribution, flag points with low likelihood as anomalous

2. Clustering approach

- Use Mahalanobis distance to take into account the covariance of the attributes
- Calculate distance of each point to the centroid, flag points with largest distance

$$d_{MH}(i, j) = \sqrt{(x(i) - x(j))^T \Sigma^{-1} (x(i) - x(j))}$$

Example

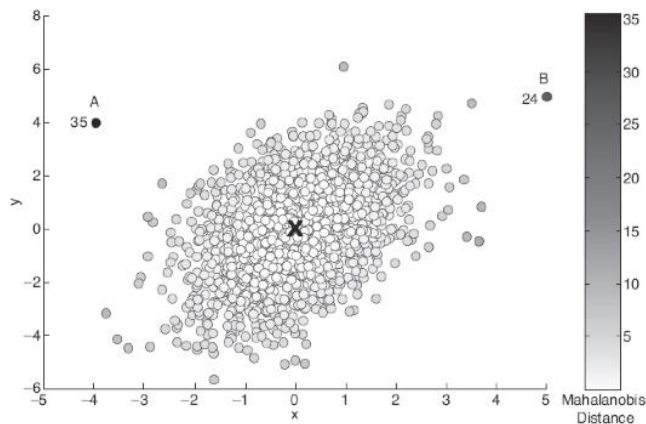
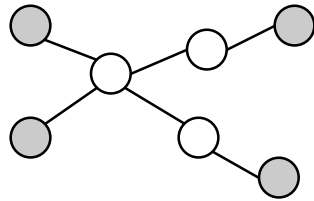


Figure 10.3. Mahalanobis distance of points from the center of a two-dimensional set of 2002 points.

Example: Network traffic (Lakhina et. al '04)



Backbone network

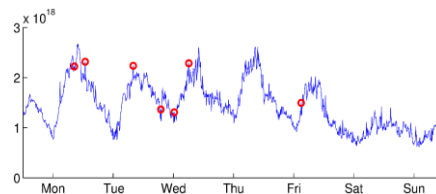
Goal: Find **source-destination** pairs with high traffic (e.g., by rate, volume)

$$Y = \begin{bmatrix} \dots \\ 100 & 30 & 42 & 212 & 1729 & 13 \\ \dots \end{bmatrix}$$

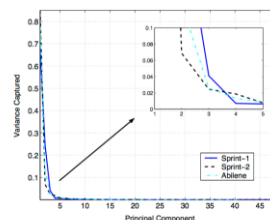
Source: Sutton, CS294 UC Berkeley

Example: Network traffic

Abilene backbone network traffic volume over 41 links collected over 4 weeks

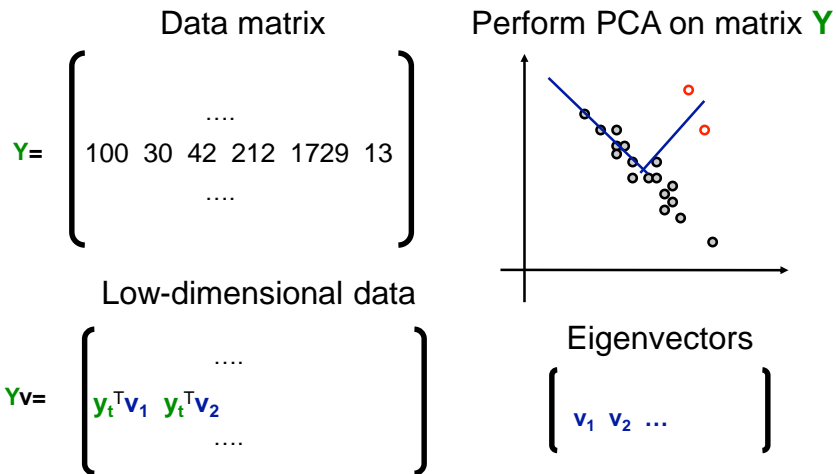


Perform PCA on 41-dim data
Select top 5 components to form "normal" subspace **P**



Source: Sutton, CS294 UC Berkeley

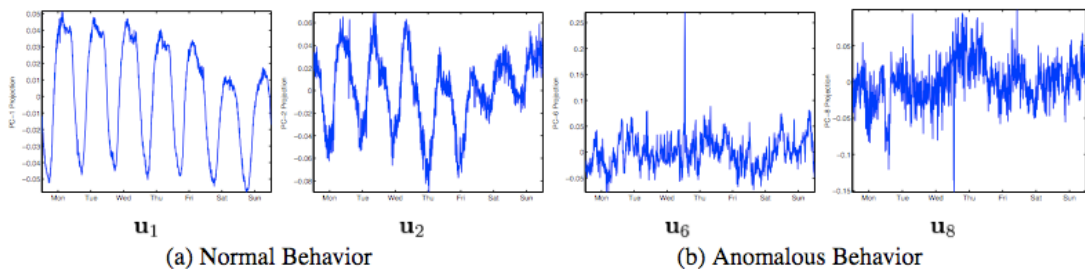
Example: Network traffic



Source: Sutton, CS294 UC Berkeley

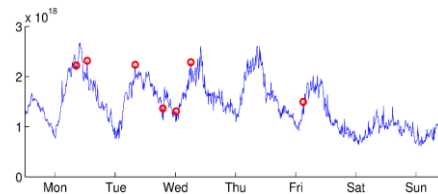
Example: Network traffic

- Projections onto principal components: $u_i = \frac{Yv_i}{\|Yv_i\|}$
 - Look for first projection that contains a 3σ from mean to identify beginning of “anomalous” subspace



Example: Network traffic

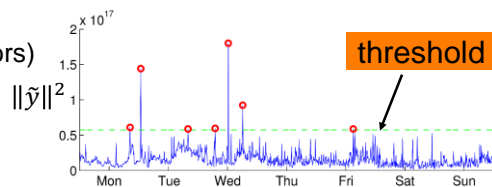
Abilene backbone network traffic volume over 41 links collected over 4 weeks



Project to residual subspace (everything but top 5 eigenvectors)

$$\hat{y} = (\mathbf{I} - \mathbf{P}\mathbf{P}^T)\mathbf{y}$$

Flag anomalies using norm of projected vector



Source: Sutton, CS294 UC Berkeley

Anomaly detection (generally)

- Knowledge representation?
 - **Categorization of data points as normal/anomalies**
- Pattern space?
 - **Set of data points** (each point is given a score)
- Score function?
 - **Distance from mean/centroid**
- Search?
 - **Exhaustive search**
consider all points, return those with distance > threshold