Outlier Analysis

Lijun Zhang <u>zlj@nju.edu.cn</u> http://cs.nju.edu.cn/zlj





Outline

Introduction

- Extreme Value Analysis
- Probabilistic Models
- Clustering for Outlier Detection
- Distance-Based Outlier Detection
- Density-Based Methods
- Information-Theoretic Models
- Outlier Validity
- Summary



Introduction (1)

A Quote

"You are unique, and if that is not fulfilled, then something has been lost."—Martha Graham

□ An Informal Definition

"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."

□ A Complementary Concept to Clustering

- Clustering attempts to determine groups of data points that are similar
- Outliers are individual data points that are different from the remaining data



Introduction (2)

Applications

- Data cleaning
 - ✓ Remove noise in data

Credit card fraud

Unusual patterns of credit card activity

Network intrusion detection

Unusual records/changes in network traffic



Introduction (3)

□ The Key Idea

- Create a model of normal patterns
- Outliers are data points that do not naturally fit within this normal model
- The "outlierness" of a data point is quantified by a outlier score
- Outputs of Outlier Detection Algorithms
 - Real-valued outlier score
 - Binary label



Outline

- Introduction
- **Extreme Value Analysis**
- Probabilistic Models
- Clustering for Outlier Detection
- Distance-Based Outlier Detection
- Density-Based Methods
- Information-Theoretic Models
- Outlier Validity
- Summary



Extreme Value Analysis (1)



- □ All extreme values are outliers
- Outliers may not be extreme values
 - {1,3,3,3,50,97,97,97,100}
 - 1 and 100 are extreme values
 - 50 is an outlier but not extreme value



Extreme Value Analysis (2)

All extreme values are outliesOutlies may not be extreme values



Univariate Extreme Value Analysis (1)



□ Statistical Tail Confidence Tests

- Suppose the density distribution is $f_X(x)$
- Tails are extreme regions s.t. $f_X(x) \le \theta$
- Symmetric Distribution
 - Two symmetric tails
 - The areas inside tails represent the cumulative probability



Univariate Extreme Value Analysis (2)



□ Statistical Tail Confidence Tests

- Suppose the density distribution is $f_X(x)$
- Tails are extreme regions s.t. $f_X(x) \le \theta$
- Asymmetric Distribution
 - Areas in two tails are different
 - Regions in the interior are not tails





The Procedure (1)

- □ A model distribution is selected
 - Normal Distribution with mean μ and standard deviation σ

$$f_X(x) = \frac{1}{\sigma \cdot \sqrt{2 \cdot \pi}} \cdot e^{\frac{-(x-\mu)^2}{2 \cdot \sigma^2}}$$

- Parameter Selection
 - Prior domain knowledge
 - Estimate from data

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} x_i$$
 $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \hat{\mu})^2$



The Procedure (2)

□ *Z*-value of a random variable

$$z_i = \frac{x_i - \mu}{\sigma}$$

- Large positive values of z_i correspond to the upper tail
- Large negative values of z_i correspond to the lower tail
- z_i follows the standard normal distribution





Multivariate Extreme Values (1

- Unimodal probability distributions with a single peak
 - Suppose the density distribution is $f_X(x)$
 - Tails are extreme regions s.t. $f_X(x) \le \theta$
- Multivariate Gaussian Distribution

$$f(\overline{X}) = \frac{1}{\sqrt{|\Sigma|} \cdot (2 \cdot \pi)^{(d/2)}} \cdot e^{-\frac{1}{2} \cdot (\overline{X} - \overline{\mu}) \Sigma^{-1} (\overline{X} - \overline{\mu})^T}$$
$$= \frac{1}{\sqrt{|\Sigma|} \cdot (2 \cdot \pi)^{(d/2)}} \cdot e^{-\frac{1}{2} \cdot Maha(\overline{X}, \overline{\mu}, \Sigma)^2}$$

where $Maha(\overline{X}, \overline{\mu}, \Sigma)$ is the Mahalanobis distance between \overline{X} and $\overline{\mu}$



Multivariate Extreme Values (2)

$\square \text{ Extreme-value Score of } \overline{X}$

- $\blacksquare Maha(\overline{X}, \overline{\mu}, \Sigma)$
- Larger values imply more extreme behavior





Multivariate Extreme Values (2

\Box Extreme-value Score of \overline{X}

- $\blacksquare Maha(\overline{X}, \overline{\mu}, \Sigma)$
- Larger values imply more extreme behavior

$\square \text{ Extreme-value Probability of } \overline{X}$

• Let \mathcal{R} be the region

 $\mathcal{R} = \{ \overline{Y} | Maha(\overline{Y}, \overline{\mu}, \Sigma) \ge Maha(\overline{X}, \overline{\mu}, \Sigma) \}$

- Cumulative probability of \mathcal{R}
- Cumulative Probability of χ^2 distribution for which the value is larger than $Maha(\bar{X}, \bar{\mu}, \Sigma)$



Why χ^2 distribution?

- □ The Mahalanobis distance
 - Let Σ be the covariance matrix

$$Maha(\overline{Y},\overline{\mu},\Sigma) = \sqrt{(\overline{Y}-\overline{\mu})\Sigma^{-1}(\overline{Y}-\overline{\mu})^{\mathsf{T}}}$$

Projection + Normalization \checkmark Let $\Sigma = U\Lambda U^{\top} = \sum_{i=1}^{d} \sigma_i^2 \mathbf{u}_i \mathbf{u}_i^{\top}$ \checkmark Then, $\Sigma^{-1} = U\Lambda^{-1}U^{\top} = \sum_{i=1}^{d} \sigma_i^{-2} \mathbf{u}_i \mathbf{u}_i^{\top}$

$$Maha(\bar{Y},\bar{\mu},\Sigma) = \sqrt{(\bar{Y}-\bar{\mu})\left(\sum_{i=1}^{d}\sigma_{i}^{-2}\mathbf{u}_{i}\mathbf{u}_{i}^{\mathsf{T}}\right)(\bar{Y}-\bar{\mu})^{\mathsf{T}}} = \sqrt{\sum_{i=1}^{d}\left(\frac{\mathbf{u}_{i}(\bar{Y}-\bar{\mu})^{\mathsf{T}}}{\sigma_{i}}\right)^{2}}$$



Adaptive to the Shape

□ *B* is an extreme value





Depth-Based Methods

Convex Hull

The convex hull of a set C, denoted conv C, is the set of all convex combinations of points in C:

conv $C = \{\theta_1 x_1 + \dots + \theta_k x_k \mid x_i \in C, \ \theta_i \ge 0, \ i = 1, \dots, k, \ \theta_1 + \dots + \theta_k = 1\}.$





The Procedure

□ The index *k* is the outlier score

Smaller values indicate a grate tendency

Algorithm FindDepthOutliers(Data Set: \mathcal{D} , Score Threshold: r) **begin**

k = 1;

repeat

Find set S of corners of convex hull of \mathcal{D} ; Assign depth k to points in S; $\mathcal{D} = \mathcal{D} - S$; k = k + 1; until(D is empty); Report points with depth at most r as outliers; end



An Example

(a)

Peeling Layers of an Onion





Limitations

No Normalization



Many data points are indistinguishable
 The computational complexity increases significantly with dimensionality



Outline

- Introduction
- **Extreme Value Analysis**
- **Probabilistic Models**
- Clustering for Outlier Detection
- Distance-Based Outlier Detection
- Density-Based Methods
- Information-Theoretic Models
- Outlier Validity
- Summary



Probabilistic Models

- Related to Probabilistic Model-Based Clustering
- □ The Key Idea
 - Assume data is generated from a mixture-based generative model
 - Learn the parameter of the model from data
 - ✓ EM algorithm
 - Evaluate the probability of each data point being generated by the model
 - Points with low values are outliers

Mixture-based Generative Model



- □ Data was generated from a mixture of k distributions with probability distribution $G_1, ..., G_k$
- □ *G_i* represents a cluster/mixture component
- \Box Each point \overline{X} is generated as follows
 - Select a mixture component with probability $\alpha_i = P(G_i)$, i = 1, ..., k
 - Assume the r-th component is selected
 - Generate a data point from G_r



Learning Parameter from Data

□ The probability that \overline{X}_j generated by the mixture model \mathcal{M} is given by

$$f^{point}(\overline{X}_{j}|\mathcal{M}) = \sum_{i=1}^{k} P(\mathcal{G}_{i}, \overline{X}_{j}) = \sum_{i=1}^{k} P(\mathcal{G}_{i}) P(\overline{X}_{j}|\mathcal{G}_{i}) = \sum_{i=1}^{k} \alpha_{i} \cdot f^{i}(\overline{X}_{j})$$

□ The probability of the data set $\mathcal{D} = {\overline{X_1}, ..., \overline{X_n}}$ generated by \mathcal{M}

$$f^{data}(\mathcal{D}|\mathcal{M}) = \prod_{j=1}^{n} f^{point}(\overline{X_j}|\mathcal{M}).$$

 $\Box \text{ Learning parameters that maximize}$ $\mathcal{L}(\mathcal{D}|\mathcal{M}) = \log(\prod_{i=1}^{n} f^{point}(\overline{X_i}|\mathcal{M})) = \sum_{i=1}^{n} \log(\sum_{i=1}^{k} \alpha_i f^i(\overline{X_i}))$

$$(\mathcal{D}|\mathcal{M}) = \log(\prod_{j=1} f^{point}(\overline{X_j}|\mathcal{M})) = \sum_{j=1} \log(\sum_{i=1} \alpha_i f^i(\overline{X_j}))$$



Identify Outliers

$\Box \text{ Outlier Score is defined as}$ $f^{point}(\overline{X_j}|\mathcal{M}) = \sum_{i=1}^{k} P(\mathcal{G}_i, \overline{X_j}) = \sum_{i=1}^{k} P(\mathcal{G}_i) P(\overline{X_j}|\mathcal{G}_i) = \sum_{i=1}^{k} \alpha_i \cdot f^i(\overline{X_j})$





Outline

- Introduction
- Extreme Value Analysis
- Probabilistic Models
- **Clustering for Outlier Detection**
- Distance-Based Outlier Detection
- Density-Based Methods
- Information-Theoretic Models
- Outlier Validity
- Summary



Clustering for Outlier Detection

Outlier Analysis v.s. Clustering

- Clustering is about finding "crowds" of data points
- Outlier analysis is about finding data points that are far away from these crowds
- Every data point is
 - Either a member of a cluster
 - Or an outlier
- Some clustering algorithms also detect outliers
 - DBSCAN, DENCLUE



The Procedure (1)

□ A Simple Way

- 1. Cluster the data
- 2. Define the outlier score as the distance of the data point to its cluster centroid





The Procedure (2)

A Better Approach

- 1. Cluster the data
- 2. Define the outlier score as the local Mahalanobis distance
 - ✓ Suppose \overline{X} belongs to cluster r

$$Maha(\overline{X}, \overline{\mu_r}, \Sigma_r) = \sqrt{(\overline{X} - \overline{\mu_r})\Sigma_r^{-1}(\overline{X} - \overline{\mu_r})^T}.$$

- $\checkmark \overline{\mu_r}$ is the mean vector of the r-th cluster
- Σ_r is the covariance matrix of the *r*-th cluster
- Multivariate Extreme Value Analysis
 Global Mahalanobis distance



A Post-processing Step

Remove Small-Size Clusters





Outline

- Introduction
- Extreme Value Analysis
- Probabilistic Models
- Clustering for Outlier Detection
- Distance-Based Outlier Detection
- Density-Based Methods
- Information-Theoretic Models
- Outlier Validity
- Summary

Distance-Based Outlier Detection



□ An *Instance-Specific* Definition

The distance-based outlier score of an object 0 is its distance to its k-th nearest neighbor
13r



Distance-Based Outlier Detection



□ An *Instance-Specific* Definition

- The distance-based outlier score of an object 0 is its distance to its k-th nearest neighbor
- Sometimes, average distance is used
- □ High-computational Cost $O(n^2)$
 - Index structure
 - Effective when the dimensionality is low
 - Pruning tricks
 - Designed for the case that only the top-r outliers are needed

The Naïve Approach for Finding Top *r*-Outliers

1. Evaluate the $n \times n$ distance matrix



The Naïve Approach for Finding Top *r*-Outliers

1. Evaluate the $n \times n$ distance matrix



2. Find the *k*-th smallest value in each row

The Naïve Approach for Finding Top *r*-Outliers

1. Evaluate the $n \times n$ distance matrix



2. Find the *k*-th smallest value in each row 3. Choose *r* data points with largest $V_k(\cdot)$



1. Evaluate a $s \times n$ distance matrix





1. Evaluate a $s \times n$ distance matrix



2. Find the *k*-th smallest value in each row



1. Evaluate a $s \times n$ distance matrix



Find the *k*-th smallest value in each row
 Identify the *r*-th score in top *s*-rows



1. Evaluate a $s \times n$ distance matrix



- 2. Find the *k*-th smallest value in each row
- **3**. Identify the *r*-th score in top *s*-rows
- 4. Remove points with $\widehat{V_k}(\cdot) \leq L_r$

Pruning Methods—Early Termination



 L_r

When completing the empty area



Pruning Methods—Early Termination



When completing the empty area



- Update $\widehat{V_k}(\cdot)$ when more distances are known
- $\Box \text{ Stop if } \widehat{V_k}(\cdot) \leq L_r$
- \Box Update L_r if necessary

Local Distance Correction Methods



Impact of Local Variations





Local Outlier Factor (LOF)

- □ Let $V^k(\overline{X})$ be the distance of \overline{X} to its k-nearest neighbor
- □ Let $L_k(\overline{X})$ be the set of points within the *k*-nearest neighbor distance of \overline{X}
- Reachability Distance
 - $R_k(\overline{X}, \overline{Y}) = \max\{Dist(\overline{X}, \overline{Y}), V^k(\overline{Y})\}\$
 - Not symmetric between \overline{X} and \overline{Y}
 - If $Dist(\overline{X}, \overline{Y})$ is large, $R_k(\overline{X}, \overline{Y}) = Dist(\overline{X}, \overline{Y})$
 - Otherwise, $R_k(\overline{X}, \overline{Y}) = V^k(\overline{Y})$
 - ✓ Smoothed out by $V^k(\overline{Y})$, more stable



Local Outlier Factor (LOF)

 $\square \text{ Average Reachability Distance}$ $AR_k(\overline{X}) = \text{MEAN}_{\overline{Y} \in L_k(\overline{X})}R_k(\overline{X}, \overline{Y})$

Local Outlier Factor

$$LOF_k(\overline{X}) = MEAN_{\overline{Y} \in L_k(\overline{X})} \frac{AR_k(\overline{X})}{AR_k(\overline{Y})}$$

Larger for Outliers
 Close to 1 for Others
 Outlier Score
 max LOF_k(X̄)



Instance-Specific Mahalanobis Distance (1)



Define a local Mahalanobis distance for each point

Based on the covariance structure of the neighborhood of a data point

□ The Challenge

- Neighborhood of a data point is hard to define with the Euclidean distance
- Euclidean distance is biased toward capturing the circular region around that point

Instance-Specific Mahalanobis Distance (2)



An agglomerative approach for neighborhood construction

• Add \overline{X} to $L^k(\overline{X})$

Data points are iteratively added to $L^k(\overline{X})$ that have the smallest distance to $L^k(\overline{X})$

 $\operatorname{argmin}_{\bar{Y}\in\mathcal{D}}\min_{\bar{Z}\in L^{k}(\bar{X})}dist(\bar{Y}-\bar{Z})$

□ Instance-specific Mahalanobis score

 $LMaha_k(\overline{X}) = Maha(\overline{X}, \overline{\mu_k(X)}, \Sigma_k(\overline{X}))$

Outlier score $\max_{k} LMaha_{k}(\overline{X})$

Instance-Specific Mahalanobis Distance (3)



Can be applied to both cases



Relation to clustering-based approaches



Outline

- Introduction
- Extreme Value Analysis
- Probabilistic Models
- Clustering for Outlier Detection
- Distance-Based Outlier Detection
- Density-Based Methods
- □ Information-Theoretic Models
- Outlier Validity
- Summary



Density-Based Methods

- □ The Key Idea
 - Determine sparse regions in the underlying data
- Limitations
 - Cannot handle variations of density



Histogram- and Grid-Based Techniques



Histogram for 1-dimensional data

Data points that lie in bins with very low frequency are reported as outliers

https://www.mathsisfun.com/data/histograms.html



- Grid for high-dimensional data
- Challenges
 - Size of grid
 - Too local
 - Sparsity





Kernel Density Estimation

 \Box Given *n* data points $\overline{X_1}, \ldots, \overline{X_n}$

$$f(\overline{X}) = \frac{1}{n} \sum_{i=1}^{n} K(\overline{X} - \overline{X_i}).$$

• $K(\cdot)$ is a kernel function

$$K(\overline{X} - \overline{X_i}) = \left(\frac{1}{h\sqrt{2\pi}}\right)^d e^{-\frac{||\overline{X} - \overline{X_i}|}{2 \cdot h^2}} .$$

Density

The density at each data point

- Computed without including the point itself in the density computation
- Low values of the density indicate greater tendency to be an outlier



Outline

- Introduction
- **Extreme Value Analysis**
- Probabilistic Models
- Clustering for Outlier Detection
- Distance-Based Outlier Detection
- Density-Based Methods
- Information-Theoretic Models
- Outlier Validity
- Summary



Information-Theoretic Models

An Example

- The 1st One: "AB 17 times"
- C in 2nd string increases its minimum description length
- Conventional Methods
 - Fix model, then calculate the deviation
- □ Information-Theoretic Models
 - Fix the deviation, then learn the model
 - Outlier score of \overline{X} : increase of the model size when \overline{X} is present



Probabilistic Models

□ The Conventional Method

- Learn the parameters of generative model with a fixed size
- Use the fit of each data point as the outlier score
- Information-Theoretic Method
 - Fix a maximum allowed deviation (a minimum value of fit)
 - Learn the size and values of parameters
 - Increase of size is used as the outlier score



Outline

- Introduction
- **Extreme Value Analysis**
- Probabilistic Models
- Clustering for Outlier Detection
- Distance-Based Outlier Detection
- Density-Based Methods
- Information-Theoretic Models
- Outlier Validity
- Summary



Outlier Validity

Methodological Challenges

- Internal criteria are rarely used in outlier analysis
- A particular validity measure will favor an algorithm using a similar objective function criterion
- Magnified because of the small sample solution space
- External Measures
 - The known outlier labels from a synthetic data set
 - The rare class labels from a real data set

Receiver Operating Characteristic (ROC) curve



- \Box *G* is the set of outliers (ground-truth)
- An algorithm outputs a outlier score
- Given a threshold t, we denote the set of outliers by S(t)
 - True-positive rate (recall)

$$TPR(t) = Recall(t) = 100 * \frac{|\mathcal{S}(t) \cap \mathcal{G}|}{|\mathcal{G}|}$$

The false positive rate

$$FPR(t) = 100 * \frac{|\mathcal{S}(t) - \mathcal{G}|}{|\mathcal{D} - \mathcal{G}|}$$

ROC Curve

Plot TPR(t) versus FPR(t)



An Example

Algorithm	Rank of ground-truth outliers
Algorithm A	1, 5, 8, 15, 20
Algorithm B	3, 7, 11, 13, 15
Random Algorithm	17, 36, 45, 59, 66
Perfect Oracle	1, 2, 3, 4, 5





Outline

- Introduction
- **Extreme Value Analysis**
- Probabilistic Models
- Clustering for Outlier Detection
- Distance-Based Outlier Detection
- Density-Based Methods
- Information-Theoretic Models
- Outlier Validity
- □ Summary



Summary

Extreme Value Analysis Univariate, Multivariate, Depth-Based Probabilistic Models Clustering for Outlier Detection Distance-Based Outlier Detection Pruning, LOF, Instance-Specific Density-Based Methods Histogram- and Grid-Based, Kernel Density Information-Theoretic Models Outlier Validity **ROC** curve