

# Outlier Analysis

---

Lijun Zhang

[zlj@nju.edu.cn](mailto:zlj@nju.edu.cn)

<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

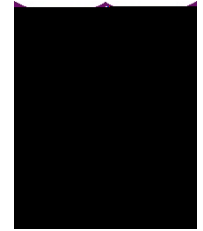
*“An outlier is an observation which deviates so much from the other observations as to appear to originate from a different mechanism.”*

## † 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)

---

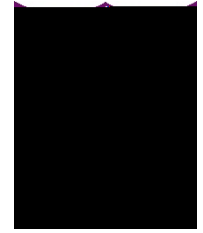


## † Applications

- „ Data cleaning
  - 9 Remove noise in data
- „ Credit card fraud
  - 9 Unusual patterns of credit card activity
- „ Network intrusion detection
  - 9 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

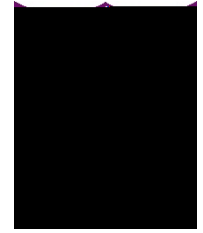
# Outline

---

- † 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)

---



## † Statistical Tails

<http://www.regentprep.org/regents/math/algtrig/ats2/normallesson.htm>

† All extreme values are outliers

† Outliers may not be extreme values

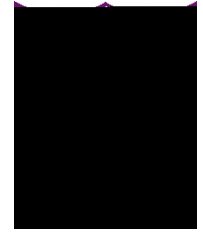
„ {1,3,3,3,50,97,97,97,100}

„ 1 and 100 are extreme values

„ 50 is an outlier but not extreme value

# Extreme Value Analysis (2)

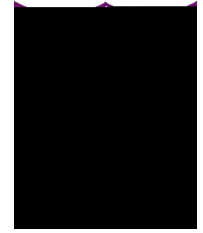
---



- † All extreme values are outliers
- † Outliers may not be extreme values



# Univariate Extreme Value Analysis (1)

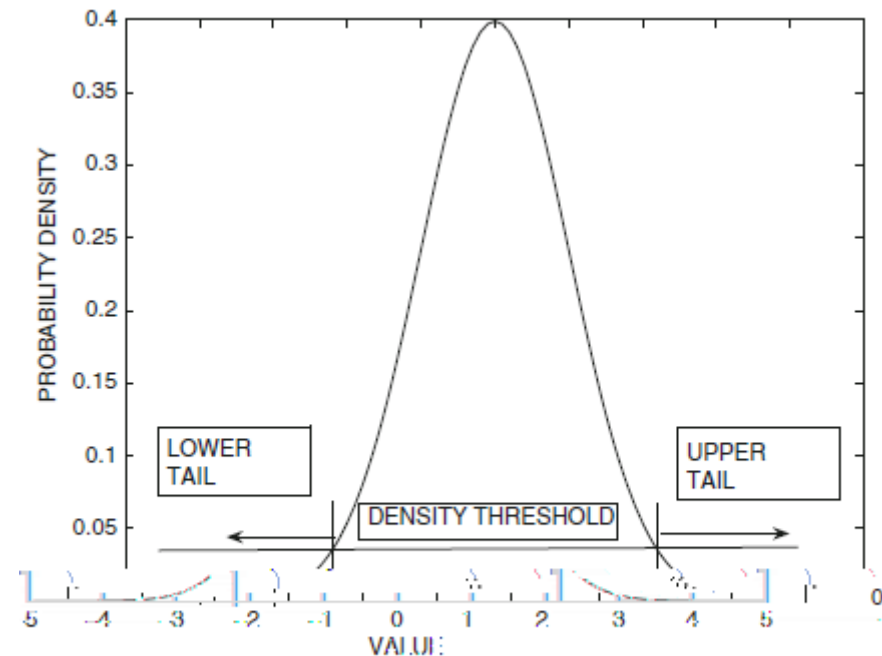


## † Statistical Tail Confidence Tests

- „ Suppose the density distribution is  $f_{\tilde{N}}(x)$
- „ Tails are **extreme** regions s.t.  $f_{\tilde{N}}(x) \leq \theta$

## † Symmetric Distribution

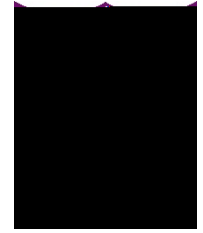
- „ Two symmetric tails
- „ The areas inside tails represent the cumulative probability



(a) Symmetric distribution

# Univariate Extreme Value Analysis (2)

---



## † Statistical Tail Confidence Tests

”

# The Procedure (1)



## † A model distribution is selected

- „ Normal Distribution with mean  $\mu$  and standard deviation  $\sigma$

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

## † Parameter Selection

- „ Prior domain knowledge
- „ Estimate from data

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

# The Procedure (2)

---

† Z-value of a random variable

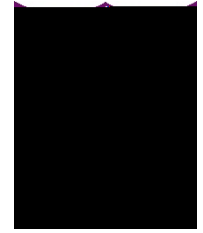
$$V_j = \frac{T_j - F}{\hat{\sigma}}$$

- „ Large positive values of  $V_j$  correspond to the upper tail
- „ Large negative values of  $V_j$  correspond to the lower tail
- „  $V_j$  follows the normal distribution

† Extreme values

- „  $|V_j| > R$

# Multivariate Extreme Values (1)



† Unimodal probability distributions with a single peak

„ Suppose the density distribution is  $f_{\tilde{N}}(x)$

„ Tails are **extreme** regions s.t.  $f_{\tilde{N}}(x) \leq \theta$

† Multivariate Gaussian Distribution

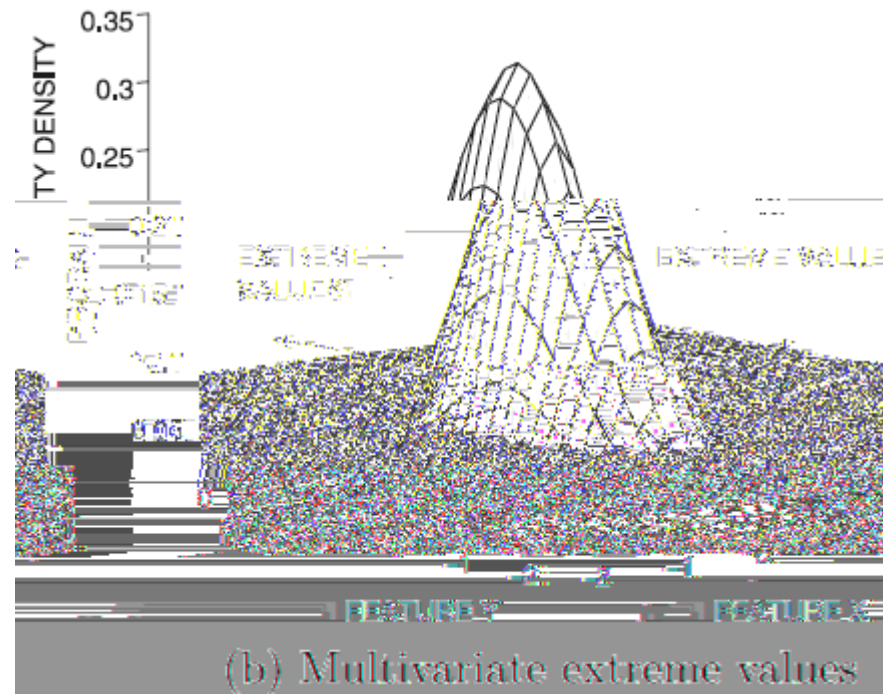
$$f(\bar{X}) = \frac{1}{\sqrt{|\Sigma|} \cdot (2 \cdot \pi)^{(d/2)}} \cdot e^{-\frac{1}{2} \cdot (\bar{X} - \bar{\mu}) \Sigma^{-1} (\bar{X} - \bar{\mu})^T}$$
$$= \frac{1}{\sqrt{|\Sigma|} \cdot (2 \cdot \pi)^{(d/2)}} \cdot e^{-\frac{1}{2} \cdot Maha(\bar{X}, \bar{\mu}, \Sigma)^2}$$

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

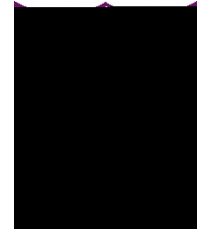
# Multivariate Extreme Values (2)

## † Extreme-value Score of $\bar{X}$

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



# Multivariate Extreme Values (2)



## † Extreme-value Score of $\bar{X}$

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

## † Extreme-value Probability of $\bar{X}$

- „ Let  $\mathcal{R}$  be the region  
$$\{ \bar{X} \in \mathbb{R}^d : Maha(\bar{X}, \bar{\mu}, \Sigma) > c \}$$
- „ Cumulative probability of  $\mathcal{R}$
- „ Cumulative Probability of  $\chi^2$  distribution for which the value is larger than  $Maha(\bar{X}, \bar{\mu}, \Sigma)$

# Why $\chi^2$ distribution?

## † The Mahalanobis distance

„ Let  $\Sigma$  be the covariance matrix

$$I = D(\hat{\mu}, \hat{\Sigma}) = \sqrt{(\hat{x} - \hat{\mu})^T \hat{\Sigma}^{-1} (\hat{x} - \hat{\mu})}$$

„ Projection+Normalization

9 Let  $L = L_1 \dots L_p$  &  $\hat{\Sigma} = \hat{\Sigma}_1 \dots \hat{\Sigma}_p$

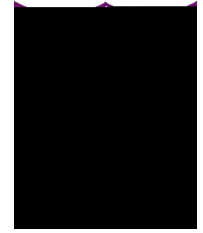
9 Then,  $L^{-1} \hat{\Sigma}^{-1} L = \hat{\Sigma}_1^{-1} \dots \hat{\Sigma}_p^{-1}$

$$I = D(\hat{\mu}, \hat{\Sigma}) = \sqrt{(\hat{x} - \hat{\mu})^T L \begin{pmatrix} \times & & \\ & \hat{\Sigma}_1^{-1} & \\ & & \ddots & \\ & & & \hat{\Sigma}_p^{-1} \end{pmatrix} (\hat{x} - \hat{\mu})} = \sqrt{\sum_{i=1}^p \left( \frac{(\hat{x} - \hat{\mu})^T L e_i}{\hat{\Sigma}_i} \right)^2}$$



# Adaptive to the Shape

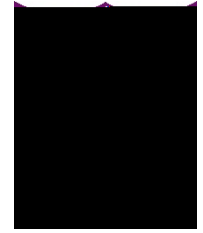
---



†  $B$  is an extreme value

# Depth-Based Methods

---



## † Convex Hull

The *convex hull* of a set  $C$ , denoted  $\text{conv } C$ , is the set of all convex combinations of points in  $C$ :

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

„ Corners

# The Procedure

---

† The index  $k$  is the outlier score

„ Smaller values indicate a greater 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

---



† Peeling Layers of an Onion

# Limitations



---

- † No Normalization
- † Many data points are indistinguishable
- † The computational complexity increases significantly with dimensionality

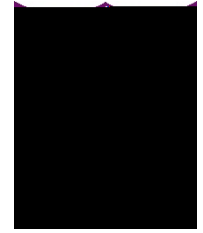
# Outline

---

- † 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

---



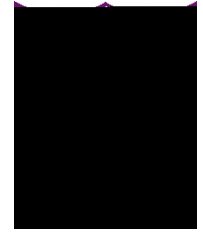
† Related to Probabilistic Model-Based Clustering

† The Key Idea

- „ Assume data is generated from a mixture-based generative model
- „ Learn the parameter of the model from data
  - 9 EM algorithm
- „ Evaluate the probability of each data point being generated by the model
  - 9 Points with low values are outliers

# Mixture-based Generative Model

---



- † Data was generated from a mixture of  $k$  distributions with probability distribution  $\mathcal{G}_1, \dots, \mathcal{G}_k$
- †  $\mathcal{G}_i$  represents a cluster/mixture component
- † Each point  $\bar{X}$  is generated as follows
  - „ Select a mixture component with probability  $\alpha_i = P(\mathcal{G}_i), i = 1, \dots, k$
  - „ Assume the  $r$ -th component is selected
  - „ Generate a data point from  $\mathcal{G}_r$



# Learning Parameter from Data

- † The probability that  $\bar{X}_j$  generated by the mixture model  $\mathcal{M}$  is given by

$$f^{point}(\bar{X}_j | \mathcal{M}) = \sum_{i=1}^k \alpha_i f^i(\bar{X}_j)$$

- † The probability of the data set  $\mathcal{D} = \{\bar{X}_1, \dots, \bar{X}_n\}$  generated by  $\mathcal{M}$

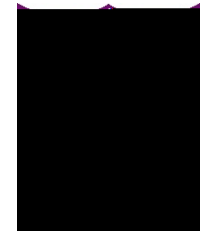
$$f^{data}(\mathcal{D} | \mathcal{M}) = \prod_{j=1}^n f^{point}(\bar{X}_j | \mathcal{M}).$$

- † Learning parameters that maximize

$$\mathcal{L}(\mathcal{D} | \mathcal{M}) = \log\left(\prod_{j=1}^n f^{point}(\bar{X}_j | \mathcal{M})\right) = \sum_{j=1}^n \log\left(\sum_{i=1}^k \alpha_i f^i(\bar{X}_j)\right)$$

# Identify Outliers

---



† Outlier Score is defined as

$$B^{\hat{a}} \hat{U}(\hat{\rho} \hat{\zeta}) L \hat{I} 2$$

$\hat{\rho} : \hat{Y}$

$:\hat{\rho}$

$s, fv \quad a, OE \quad a \pm t '5$

# Outline

---

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

# 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

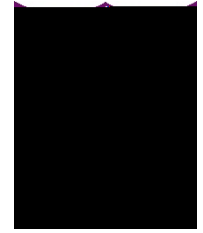
- „ Either a member of a cluster
- „ Or an outlier

## † 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**

9 Suppose  $x$  belongs to cluster  $N$

$$Maha(\bar{X}, \bar{\mu}_r, \Sigma_r) = \sqrt{(\bar{X} - \bar{\mu}_r) \Sigma_r^{-1} (\bar{X} - \bar{\mu}_r)^T}$$

9  $\bar{\mu}_a$  is the mean vector of the  $N$ th cluster

9  $\Sigma_a$  is the covariance matrix of the  $N$ th cluster

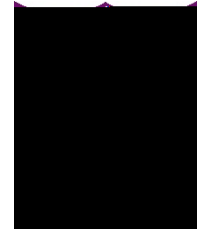
## † Multivariate Extreme Value Analysis

„ **Global Mahalanobis distance**

# A Post-processing Step

---

† Remove Small-Size Clusters



# Outline

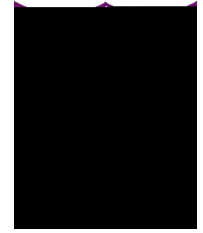
---

- † 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

---



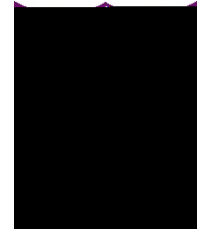
## † An Instance-Specific Definition

- „ The distance-based outlier score of an object  $O$  is its distance to its  $k$ -th nearest neighbor

$$k > 3$$

# Distance-Based Outlier Detection

---



## † An Instance-Specific Definition

- „ The distance-based outlier score of an object  $O$  is its distance to its  $k$ -th nearest neighbor
- „ Sometimes, average distance is used

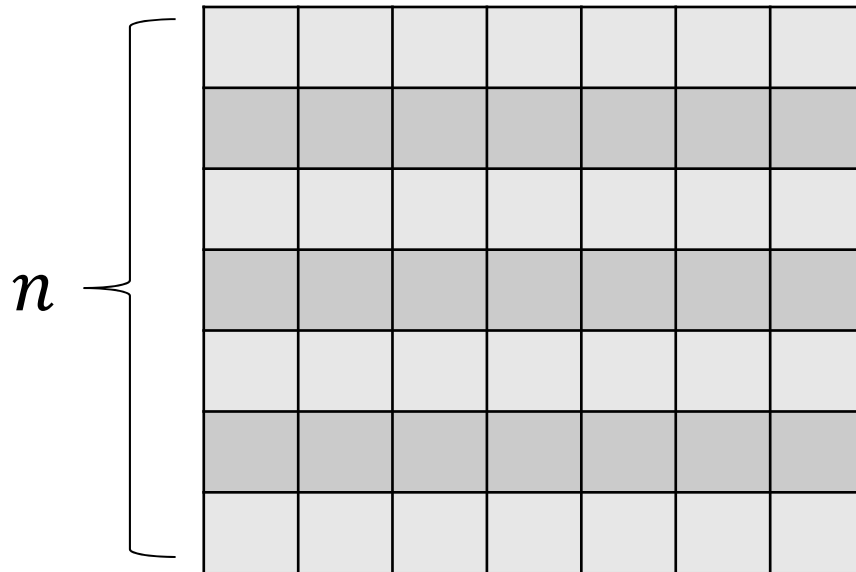
## † High-computational Cost $O(n^6)$

- „ Index structure
  - 9 Effective when the dimensionality is low
- „ Pruning tricks
  - 9 Designed for the case that only the top-outliers are needed

# The Naïve Approach for Finding Top $r$ -Outliers



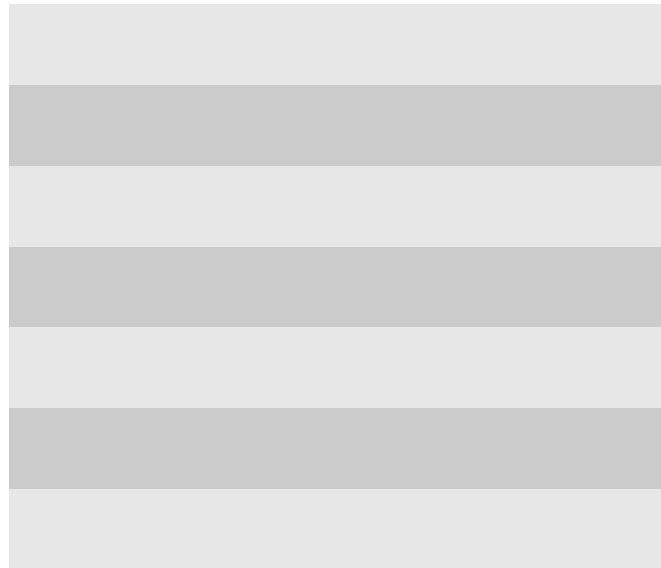
1. Evaluate the  $n \times n$  distance matrix



# The Naïve Approach for Finding Top $r$ -Outliers

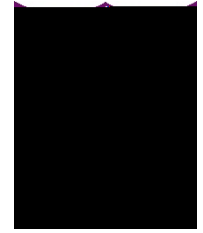


1. Evaluate the  $n \times n$  distance matrix

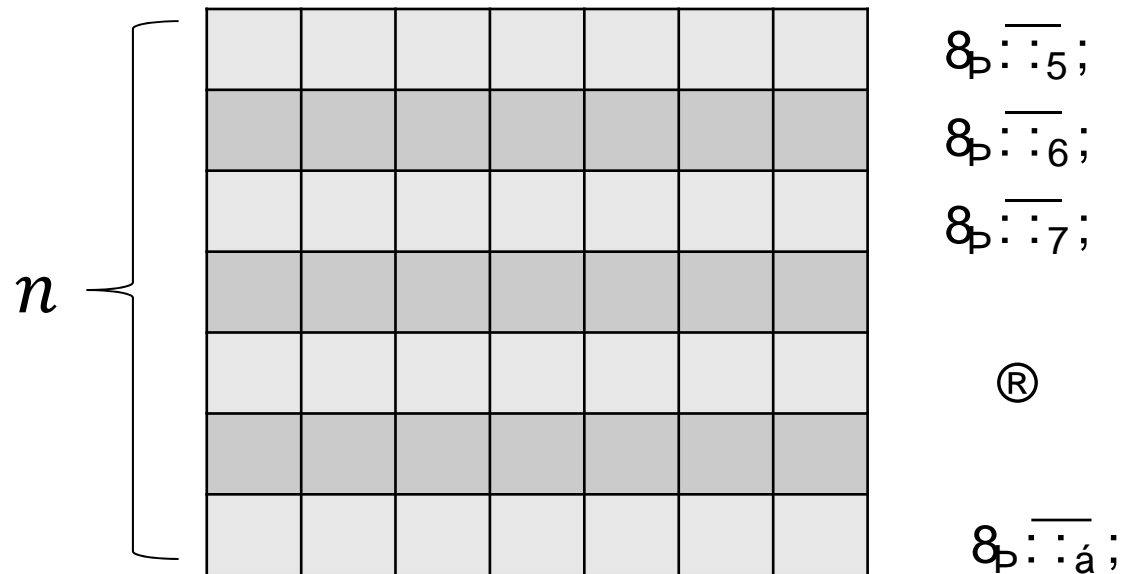


2. Find the  $k$ -th **smallest** value in each row

# The Naïve Approach for Finding Top $r$ -Outliers



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

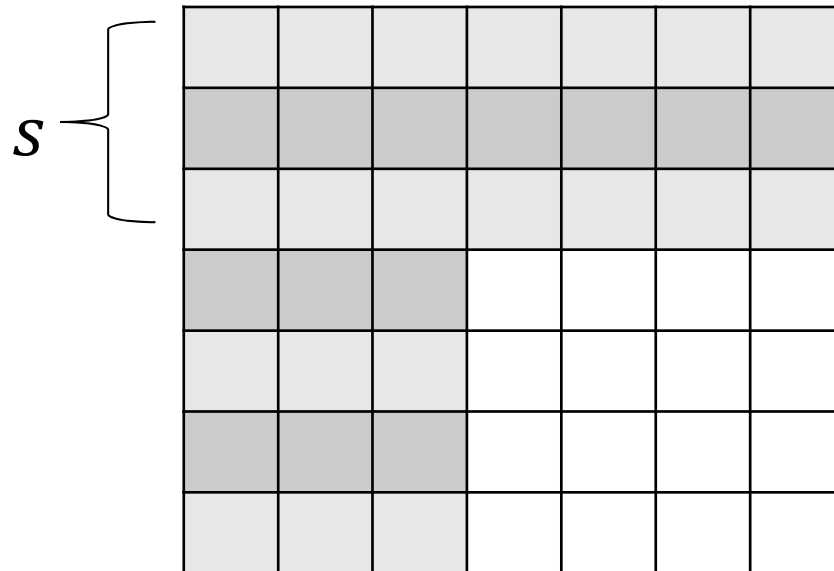


2. Find the  $k$ -th **smallest** value in each row
3. Choose  $r$  data points with **largest**  $V_{\mathcal{P}}(\cdot)$

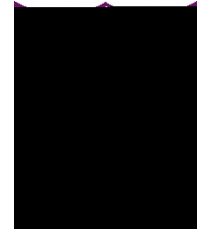
# Pruning Methods—Sampling



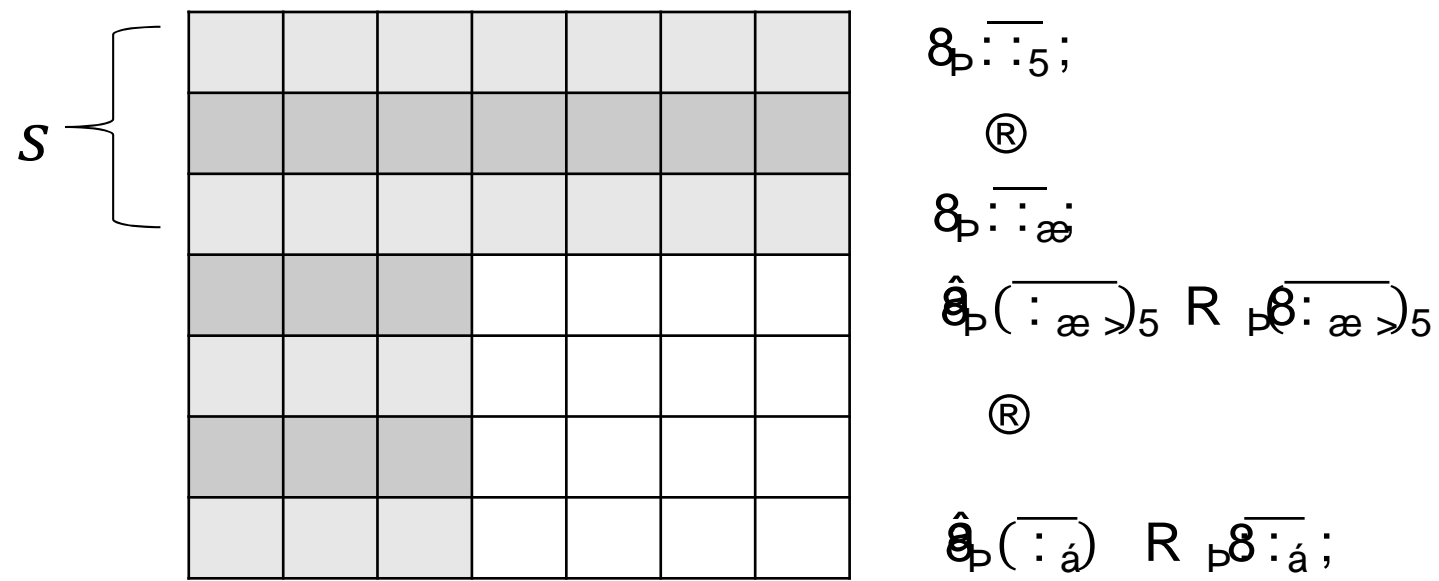
1. Evaluate a  $s \times n$  distance matrix



# Pruning Methods—Sampling



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

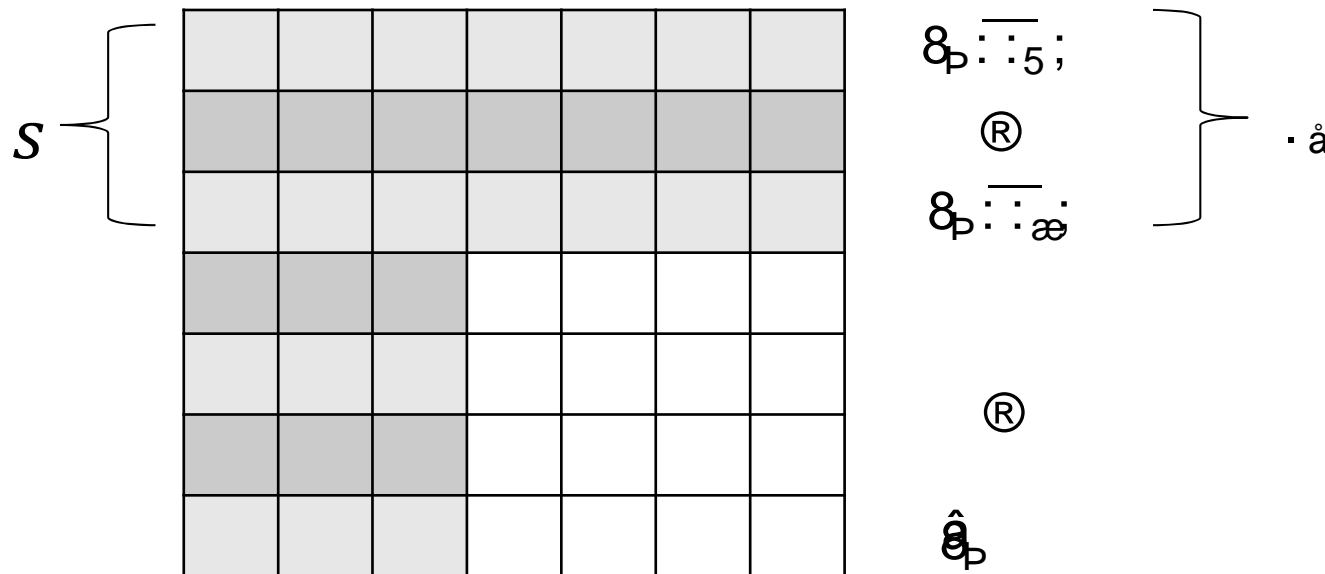


2. Find the  $k$ -th smallest value in each row

# Pruning Methods—Sampling



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

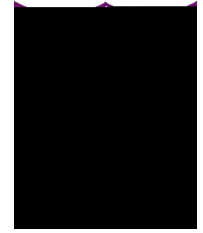


2. Find the  $k$ -th smallest value in each row

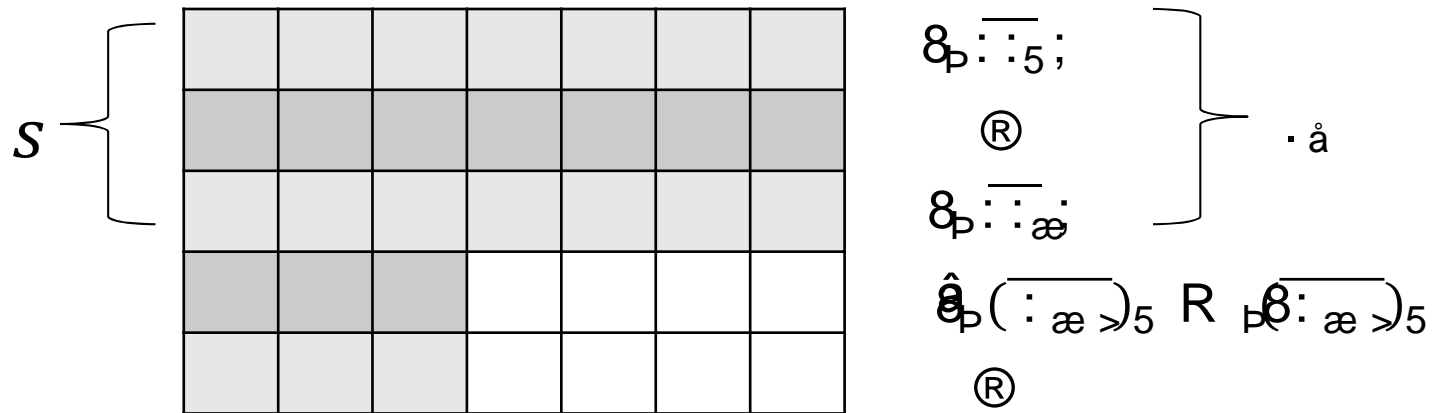
3. Identify the  $r$ -th score in top  $s$ -rows



# Pruning Methods—Sampling



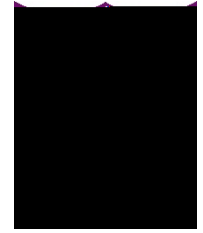
1. Evaluate a  $s \times n$  distance matrix



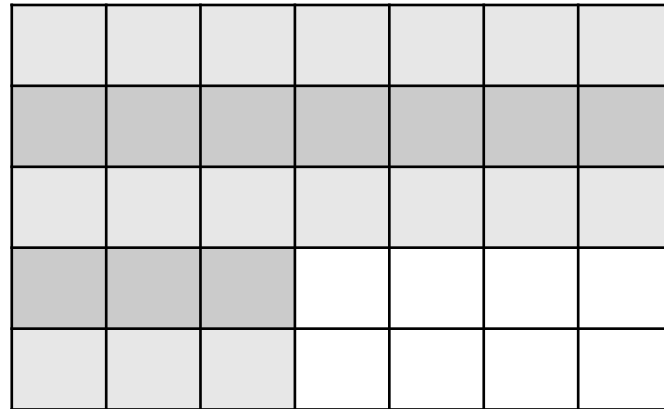
2. Find the  $k$ -th smallest value in each row
3. Identify the  $r$ -th score in top  $s$ -rows
4. Remove points with  $\hat{V}_P(\cdot) \leq L_a$

# Pruning Methods—Early Termination

---



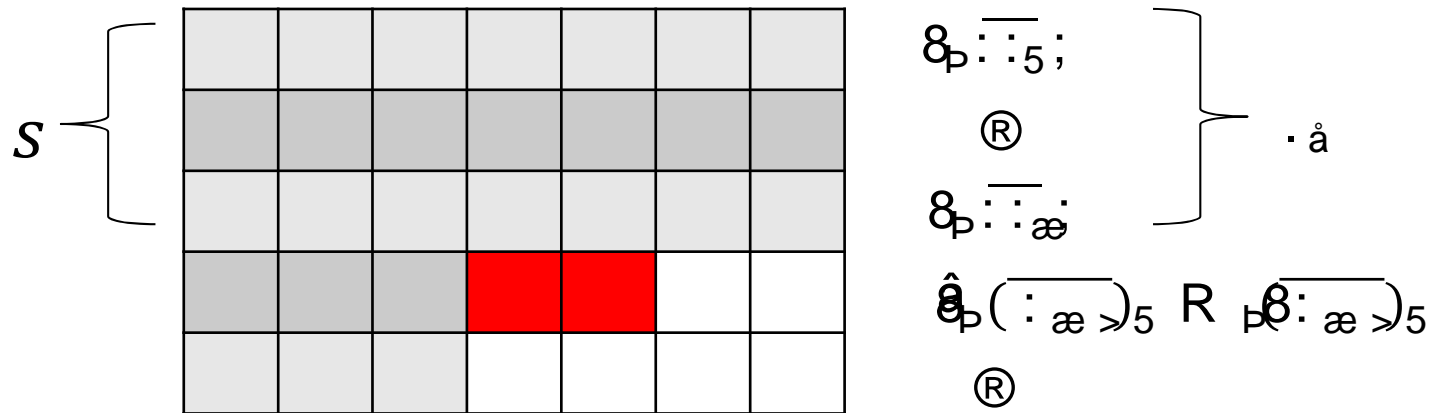
† When completing the empty area



# Pruning Methods—Early Termination



† When completing the empty area



† Update  $\widehat{V}_p(\cdot)$  when more distances are known

† Stop if  $\widehat{V}_p(\cdot) \leq L_{\dot{a}}$

† Update  $L_{\dot{a}}$  if necessary

# Local Distance Correction Methods

---



† Impact of Local Variations

# Local Outlier Factor (LOF)

- † Let  $V^p(\bar{X})$  be the distance of  $\bar{X}$  to its  $k$ -nearest neighbor
- † Let  $L_p(\bar{X})$  be the set of points within the  $k$ -nearest neighbor distance of  $\bar{X}$
- † Reachability Distance

$$R_p(\bar{X}, \bar{Y}) = \max\{Dist(\bar{X}, \bar{Y}), V^k(\bar{Y})\}$$

- „ Not symmetric between  $\bar{X}$  and  $\bar{Y}$
- „ If  $Dist(\bar{X}, \bar{Y})$  is large,  $R_p(\bar{X}, \bar{Y}) = Dist(\bar{X}, \bar{Y})$
- „ Otherwise,  $R_p(\bar{X}, \bar{Y}) = V^p(\bar{Y})$
- 9 Smoothed out by  $\delta^p$ : \$, more stable

# Local Outlier Factor (LOF)

---

## † Average Reachability Distance

$$AR_k(\bar{X}) = \text{MEAN}_{\bar{Y} \in L_k(\bar{X})} R_k(\bar{X}, \bar{Y})$$

## † Local Outlier Factor

$$LOF_k(\bar{X}) = \text{MEAN}_{\bar{Y} \in L_k(\bar{X})} \frac{AR_k(\bar{X})}{AR_k(\bar{Y})}$$

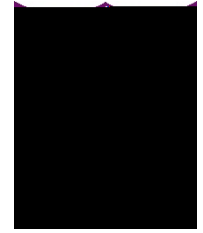
- „ Larger for Outliers
- „ Close to 1 for Others

## † Outlier Score

$$\max_{\bar{p}} LOF_{\bar{p}}(\bar{X})$$

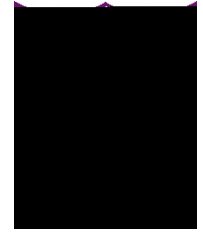
# Instance-Specific Mahalanobis Distance (1)

---



- † 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  $\bar{X}$  to  $L^P(\bar{X})$

„ Data points are **iteratively** added to  $L^P(\bar{X})$  that have the smallest distance to  $L^P(\bar{X})$



† Instance-specific Mahalanobis score

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

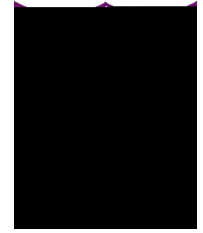
† Outlier score

$$\bullet \int_P \check{.} / = D_{\mathbb{F}} :: \$;$$



# 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

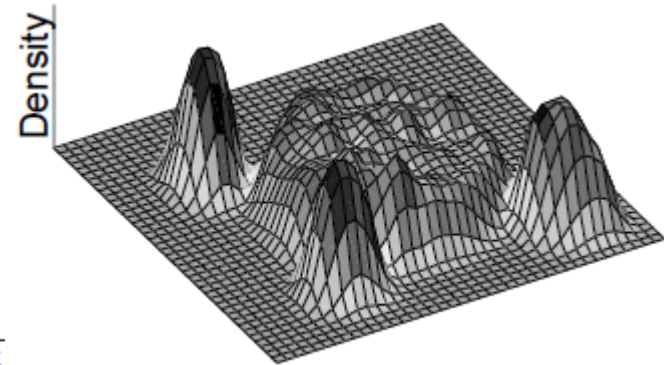
# Kernel Density Estimation

† Given  $n$  data points  $\bar{X}_1, \dots, \bar{X}_n$

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

„  $K(\cdot)$  is a kernel function

$$K(\bar{X} - \bar{X}_i) = \left( \frac{1}{h\sqrt{2\pi}} \right)^d e^{-\frac{\|\bar{X} - \bar{X}_i\|^2}{2h^2}}$$



† 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

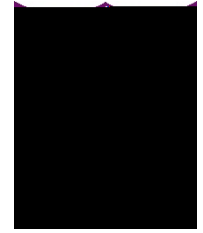
# Outline

---

- † 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

ABABABABABABABABABABABABABABABABABAB  
ABABACABABABABABABABABABABABABABABAB

- „ The 1<sup>st</sup> One: “AB 17 times”
- „ C in 2<sup>nd</sup> string increases its minimum description length

## † Conventional Methods

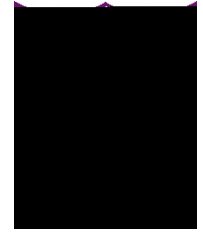
- „ Fix model, then calculate the deviation

## † Information-Theoretic Models

- „ Fix the deviation, then learn the model
- „ Outlier score of  $\Delta$  increase of the model size when  $\Delta$  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



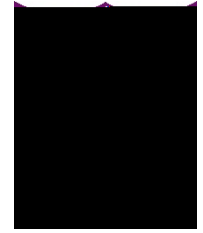
# Outline

---

- † 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



- †  $\mathcal{G}$  is the set of outliers (ground-truth)
- † An algorithm outputs an outlier score
- † Given a threshold  $t$ , we denote the set of outliers by  $\mathcal{S}(t)$

- „ 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  $FPR$  versus  $TPR$

# 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