

Ludwig-Maximilians-Universität München
Lehrstuhl für Datenbanksysteme und Data Mining
Prof. Dr. Thomas Seidl

Knowledge Discovery and Data Mining I

Winter Semester 2018/19

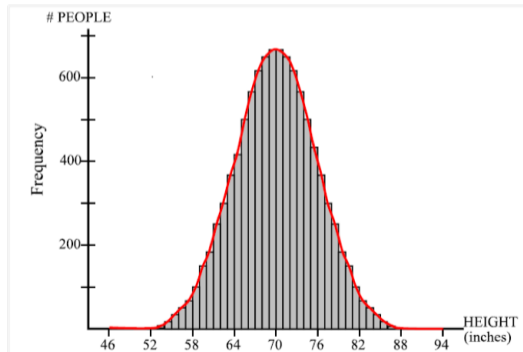


Introduction

What is an outlier?

Hawkins (1980) "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."

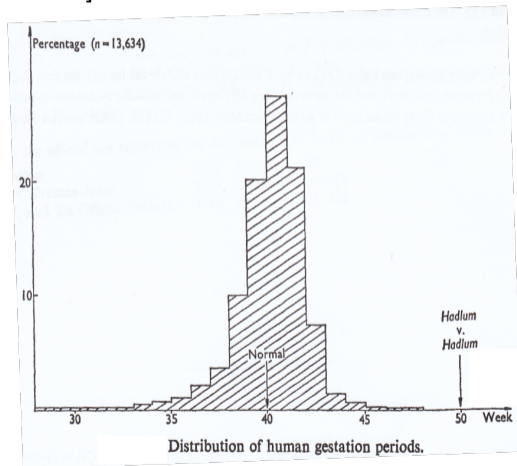
- ▶ Statistics-based intuition:
 - ▶ Normal data objects follow a "generating mechanism", e.g. some given statistical process
 - ▶ Abnormal objects deviate from this generating mechanism



Introduction

Example: Hadlum vs. Hadlum (1949) [Barnett 1978]

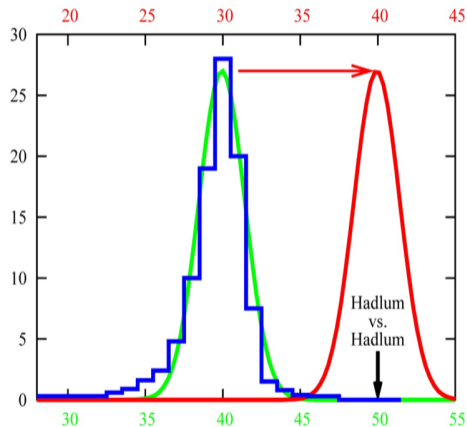
- ▶ The birth of a child to Mrs. Hadlum happened 349 days after Mr. Hadlum left for military service.
- ▶ Average human gestation period is 280 days (40 weeks).
- ▶ Statistically, 349 days is an outlier.



Introduction

Example: Hadlum vs. Hadlum (1949) [Barnett 1978]

- ▶ Blue: statistical basis (13634 observations of gestation periods)
- ▶ Green: assumed underlying Gaussian process
 - ▶ Very low probability for the birth of Mrs. Hadlums child being generated by this process
- ▶ Red: assumption of Mr. Hadlum (another Gaussian process responsible for the observed birth, where the gestation period starts later)



Applications

- ▶ Fraud detection
 - ▶ Purchasing behavior of a credit card owner usually changes when the card is stolen
 - ▶ Abnormal buying patterns can characterize credit card abuse
- ▶ Medicine
 - ▶ Whether a particular test result is abnormal may depend on other characteristics of the patients (e.g. gender, age, ...)
 - ▶ Unusual symptoms or test results may indicate potential health problems of a patient
- ▶ Public health
 - ▶ The occurrence of a particular disease, e.g. tetanus, scattered across various hospitals of a city indicate problems with the corresponding vaccination program in that city
 - ▶ Whether an occurrence is abnormal depends on different aspects like frequency, spatial correlation, etc.

Applications (cont'd)

- ▶ Sports statistics
 - ▶ In many sports, various parameters are recorded for players in order to evaluate the players' performances
 - ▶ Outstanding (in a positive as well as a negative sense) players may be identified as having abnormal parameter values
 - ▶ Sometimes, players show abnormal values only on a subset or a special combination of the recorded parameters
- ▶ Detecting measurement errors
 - ▶ Data derived from sensors (e.g. in a given scientific experiment) may contain measurement errors
 - ▶ Abnormal values could provide an indication of a measurement error
 - ▶ Removing such errors can be important in other data mining and data analysis tasks
 - ▶ *"One person's noise could be another person's signal."*

Important Properties of Outlier Models

- ▶ Global vs. local approach
 - ▶ "Outlierness" regarding whole dataset (global) or regarding a subset of data (local)?
- ▶ Labeling vs. Scoring
 - ▶ Binary decision or outlier degree score?
- ▶ Assumptions about "Outlierness"
 - ▶ What are the characteristics of an outlier object?

Agenda

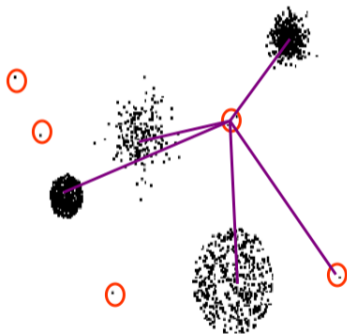
1. Introduction
2. Basics
3. Unsupervised Methods
 - 3.1 Frequent Pattern Mining
 - 3.2 Clustering
 - 3.3 **Outlier Detection**
 - 3.3.1 **Clustering-based Outliers**
 - 3.3.2 Statistical Outliers
 - 3.3.3 Distance-based Outliers
 - 3.3.4 Density-based Outliers
 - 3.3.5 Angle-based Outliers
 - 3.3.6 Summary
4. Supervised Methods
5. Advanced Topics

Clustering-based Outliers

An object is a cluster-based outlier if it does not strongly belong to any cluster.

Basic Idea

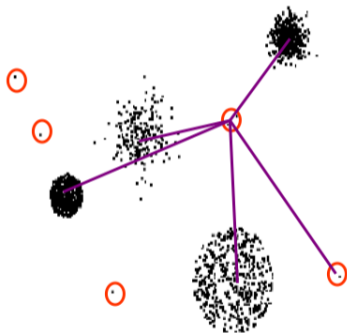
- ▶ Cluster the data into groups
- ▶ Choose points in small clusters as candidate outliers.
- ▶ Compute the distance between candidate points and non-candidate clusters.
- ▶ If candidate points are far from all other non-candidate points and clusters, they are outliers



Clustering-based Outliers

More Systematic Approaches

- ▶ Find clusters and then assess the degree to which a point belongs to any cluster
 - ▶ E.g. for k -Means, use distance to the centroid
- ▶ If eliminating a point results in substantial improvement of the objective function, we could classify it as an outlier
 - ▶ Clustering creates a model of the data and the outliers distort that model.
 - ▶ Applicable to clustering algorithms optimizing some objective function (e.g. k -means)



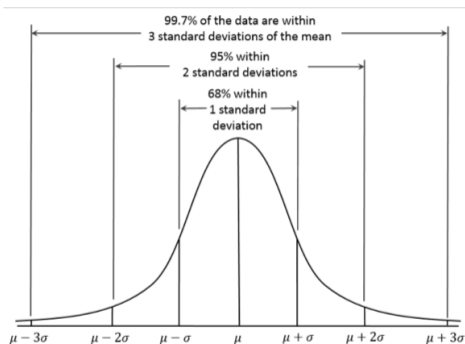
Agenda

1. Introduction
2. Basics
3. Unsupervised Methods
 - 3.1 Frequent Pattern Mining
 - 3.2 Clustering
 - 3.3 **Outlier Detection**
 - 3.3.1 Clustering-based Outliers
 - 3.3.2 **Statistical Outliers**
 - 3.3.3 Distance-based Outliers
 - 3.3.4 Density-based Outliers
 - 3.3.5 Angle-based Outliers
 - 3.3.6 Summary
4. Supervised Methods
5. Advanced Topics

Statistical Tests

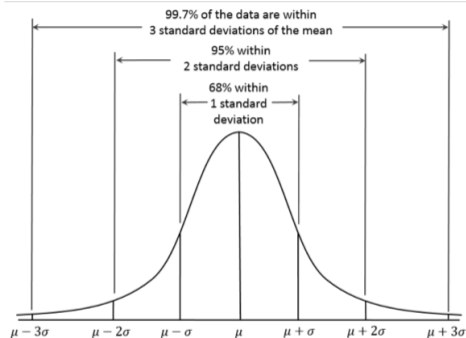
General Idea

- ▶ Given a certain kind of statistical distribution (e.g., Gaussian)
- ▶ Compute the parameters assuming all data points have been generated by such a statistical distribution (e.g., mean and standard deviation)
- ▶ Outliers are points that have a low probability to be generated by the overall distribution (e.g., deviate more than 3 times the standard deviation from the mean)



Basic Assumption

- ▶ Normal data objects follow a (known) distribution and occur in a high probability region of this model
- ▶ Outliers deviate strongly from this distribution



Statistical Tests

A huge number of different tests are available differing in

- ▶ Type of data distribution (e.g. Gaussian)
- ▶ Number of variables, i.e., dimensions of the data objects (univariate/multivariate)
- ▶ Number of distributions (mixture models)
- ▶ Parametric versus non-parametric (e.g. histogram-based)

Example on the Following Slides

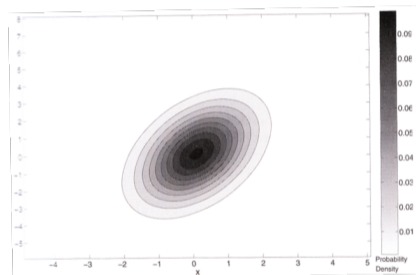
- ▶ Gaussian distribution
- ▶ Multivariate
- ▶ Single model
- ▶ Parametric

Statistical Outliers: Gaussian Distribution

Probability Density Function of a Multivariate Normal Distribution

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

- ▶ μ is the mean value of all points (usually data are normalized such that $\mu = 0$)
- ▶ Σ is the covariance matrix from the mean



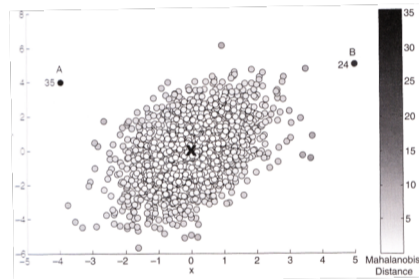
Statistical Outliers: Mahalanobis Distance

Mahalanobis Distance

Mahalanobis distance of point x to μ :

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

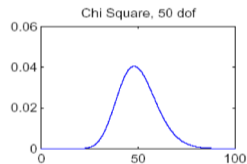
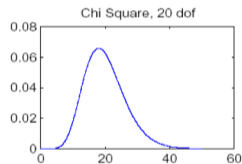
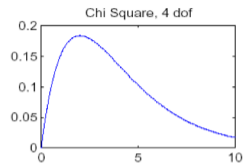
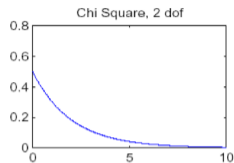
- ▶ $MDist$ follows a χ^2 -distribution with d degrees of freedom ($d =$ data dimensionality)
- ▶ Outliers: All points x , with $MDist(x, \mu) > \chi^2(0.975) (\approx 3\sigma)$



Statistical Outliers: Problems

Problems

- ▶ Curse of dimensionality: The larger the degree of freedom, the more similar the $MDist$ values for all points

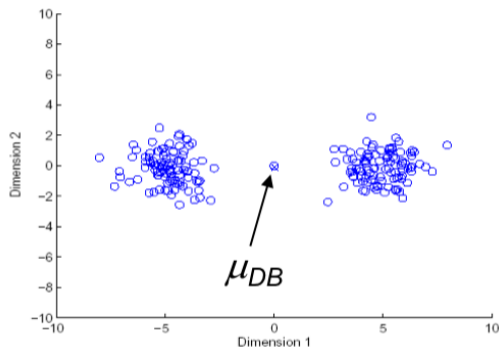


- ▶ x-axis = observed $MDist$ values
- ▶ y-axis = frequency of observation

Statistical Outliers: Problems

Problems (cont'd)

- ▶ Robustness
 - ▶ Mean and standard deviation are very sensitive to outliers
 - ▶ These values are computed for the complete data set (including potential outliers)
 - ▶ The *MDist* is used to determine outliers although the *MDist* values are influenced by these outliers



Statistical Outliers: Problems

Problems (cont'd)

- ▶ Data distribution is fixed
- ▶ Low flexibility (if no mixture models)
- ▶ Global method

Agenda

1. Introduction
2. Basics
3. Unsupervised Methods
 - 3.1 Frequent Pattern Mining
 - 3.2 Clustering
 - 3.3 **Outlier Detection**
 - 3.3.1 Clustering-based Outliers
 - 3.3.2 Statistical Outliers
 - 3.3.3 **Distance-based Outliers**
 - 3.3.4 Density-based Outliers
 - 3.3.5 Angle-based Outliers
 - 3.3.6 Summary
4. Supervised Methods
5. Advanced Topics

Distance-Based Approaches

General Idea

Judge a point based on the distance(s) to its neighbors (Several variants proposed)

Basic Assumption

- ▶ Normal data objects have a dense neighborhood
- ▶ Outliers are far apart from their neighbors, i.e., have a less dense neighborhood

Distance-Based Approaches

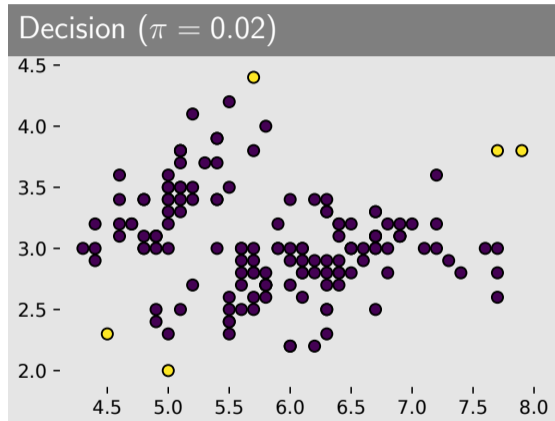
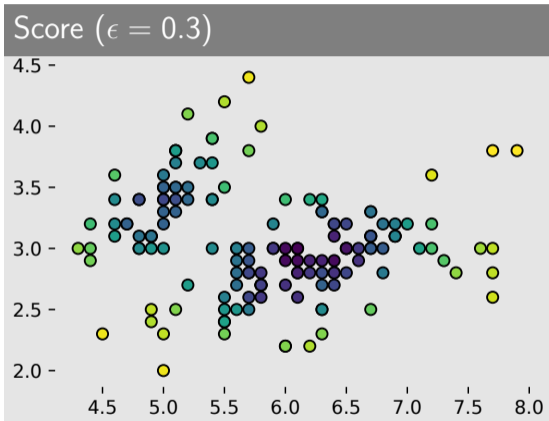
$D(\epsilon, \pi)$ Outliers²⁰

- ▶ Given: radius ϵ , percentage π
- ▶ A point p is considered an outlier if at most π percent of all points (including p) have a distance to p less than ϵ .

$$OutlierSet(\epsilon, \pi) = \left\{ p \in D \mid \frac{|\{q \in D \mid dist(p, q) < \epsilon\}|}{|D|} \leq \pi \right\}$$

²⁰E. Knorr, R. Ng. *A Unified Notion of Outliers: Properties and Computation*. KDD'97

Distance-Based Approaches: $D(\epsilon, \pi)$ Example



Distance-Based Approaches: k NN

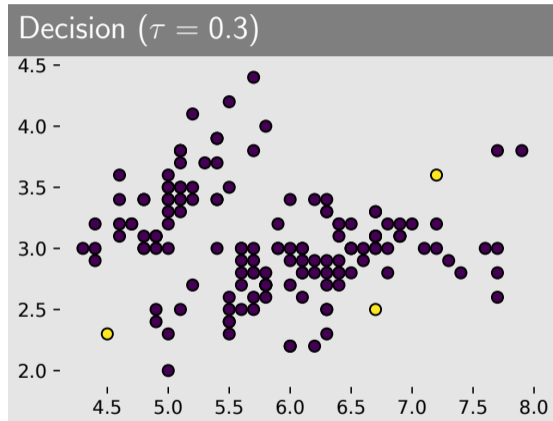
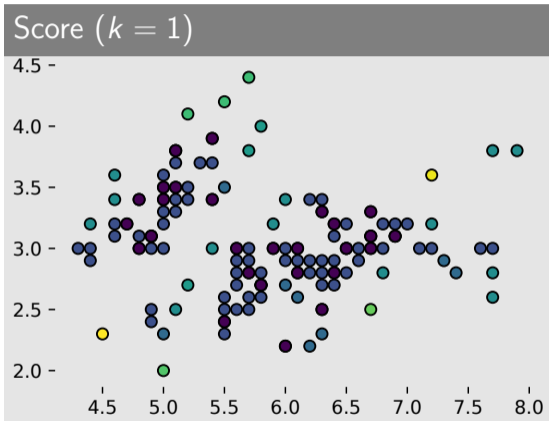
Outlier scoring based on k NN distances

General models: Take the k NN distance of a point as its outlier score

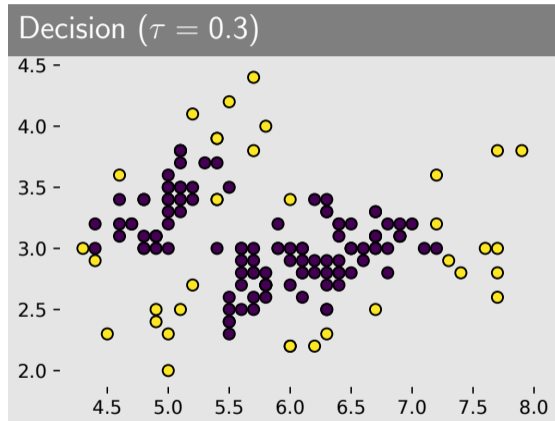
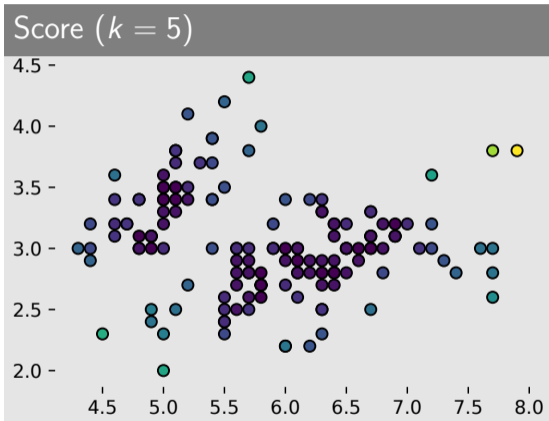
Decision

k -distance above some threshold $\tau \iff$ Outlier.

Distance-Based Approaches: k NN Example



Distance-Based Approaches: k NN Example



kNN: Problems

Problems

- ▶ Highly sensitive towards k :
 - ▶ Too small k : small number of close neighbors can cause low outlier scores.
 - ▶ Too large: all objects in a cluster with less than k objects might become outliers.
- ▶ cannot handle datasets with regions of widely different densities due to the global threshold

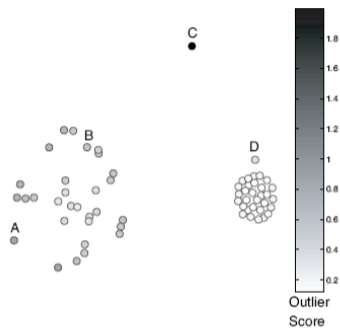


Figure 10.7. Outlier score based on the distance to the fifth nearest neighbor. Clusters of differing density.

Image Source: P. Tan, M. Steinbach, V. Kumar (2006). *Classification: basic concepts, decision trees, and model evaluation. Introduction to data mining*, 1, 145-205.

Agenda

1. Introduction
2. Basics
3. Unsupervised Methods
 - 3.1 Frequent Pattern Mining
 - 3.2 Clustering
 - 3.3 **Outlier Detection**
 - 3.3.1 Clustering-based Outliers
 - 3.3.2 Statistical Outliers
 - 3.3.3 Distance-based Outliers
 - 3.3.4 **Density-based Outliers**
 - 3.3.5 Angle-based Outliers
 - 3.3.6 Summary
4. Supervised Methods
5. Advanced Topics

Density-Based Approaches

General Idea

- ▶ Compare the density around a point with the density around its local neighbors.
- ▶ The relative density of a point compared to its neighbors is computed as an outlier score.
- ▶ Approaches also differ in how to estimate density.

Basic Assumption

- ▶ The density around a normal data object is similar to the density around its neighbors.
- ▶ The density around an outlier is considerably different to the density around its neighbors.

Density-Based Approaches

Problems

- ▶ Different definitions of density: e.g., #points within a specified distance ϵ from the given object
- ▶ The choice of ϵ is critical (too small \implies normal points considered as outliers; too big \implies outliers considered normal)
- ▶ A global notion of density is problematic (as it is in clustering); fails when data contain regions of different densities

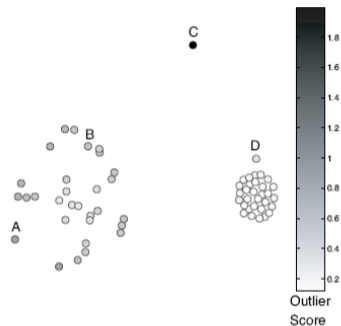


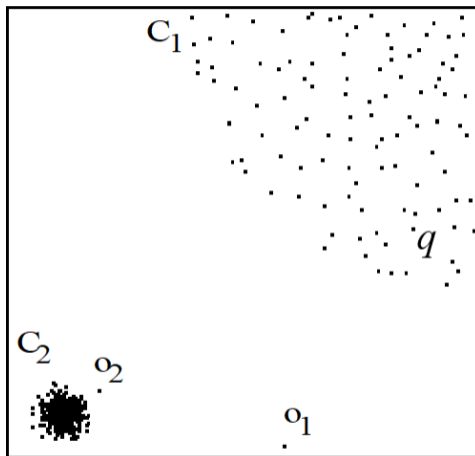
Figure 10.7. Outlier score based on the distance to the fifth nearest neighbor. Clusters of differing density.

D has a higher absolute density than *A* but compared to its neighborhood, *D*'s density is lower.

Density-Based Approaches

Failure Case of Distance-Based

- ▶ $D(\epsilon, \pi)$: parameters ϵ, π cannot be chosen s.t. o_2 is outlier, but none of the points in C_1 (e.g. q)
- ▶ k NN-distance: k NN-distance of objects in C_1 (e.g. q) larger than the k NN-distance of o_2 .



Density-Based Approaches

Solution

Consider the relative density w.r.t. to the neighbourhood.

Model

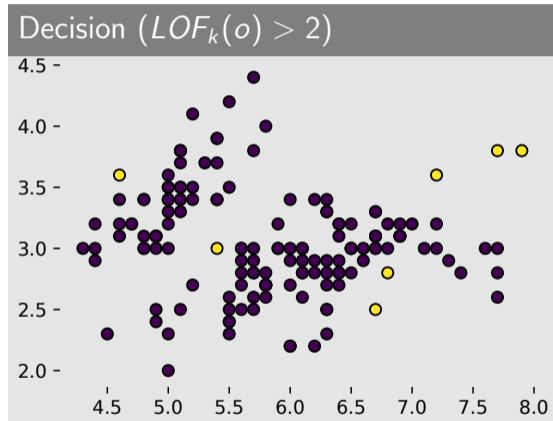
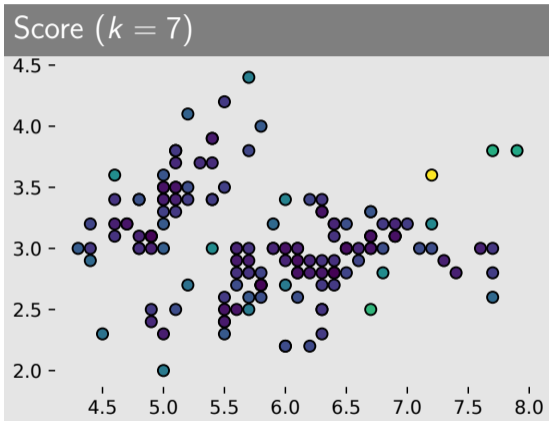
- ▶ Local Density (ld) of point p (inverse of avg. distance of k NNs of p)

$$ld_k(p) = \left(\frac{1}{k} \sum_{o \in kNN(p)} dist(p, o) \right)^{-1}$$

- ▶ Local Outlier Factor (LOF) of p (avg. ratio of ld s of k NNs of p and ld of p)

$$LOF_k(p) = \frac{1}{k} \sum_{o \in kNN(p)} \frac{ld_k(o)}{ld_k(p)}$$

Density-Based Approaches



Density-Based Approaches

Extension (Smoothing factor)

- ▶ Reachability "distance"

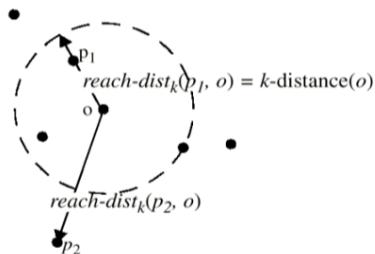
$$rd_k(p, o) = \max\{kdist(o), dist(p, o)\}$$

- ▶ Local reachability distance lrd_k

$$lrd_k(p) = \left(\frac{1}{k} \sum_{o \in kNN(p)} rd(p, o) \right)^{-1}$$

- ▶ Replace ld by lrd

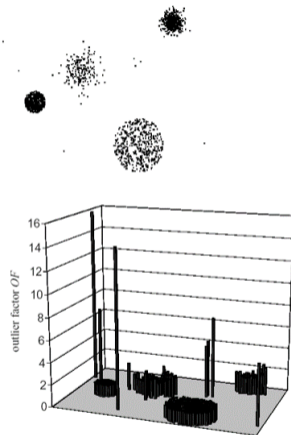
$$LOF_k(p) = \frac{1}{k} \sum_{o \in kNN(p)} \frac{lrd_k(o)}{lrd_k(p)}$$



Density-Based Approaches

Discussion

- ▶ $LOF \approx 1 \implies$ point in cluster
- ▶ $LOF \gg 1 \implies$ outlier.
- ▶ Choice of k defines the reference set



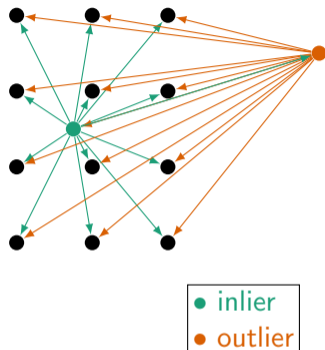
Agenda

1. Introduction
2. Basics
3. Unsupervised Methods
 - 3.1 Frequent Pattern Mining
 - 3.2 Clustering
 - 3.3 **Outlier Detection**
 - 3.3.1 Clustering-based Outliers
 - 3.3.2 Statistical Outliers
 - 3.3.3 Distance-based Outliers
 - 3.3.4 Density-based Outliers
 - 3.3.5 **Angle-based Outliers**
 - 3.3.6 Summary
4. Supervised Methods
5. Advanced Topics

Angle-Based Approach

General Idea

- ▶ Angles are more stable than distances in high dimensional spaces
- ▶ *o outlier* if most other objects are located in similar directions
- ▶ *o no outlier* if many other objects are located in varying directions



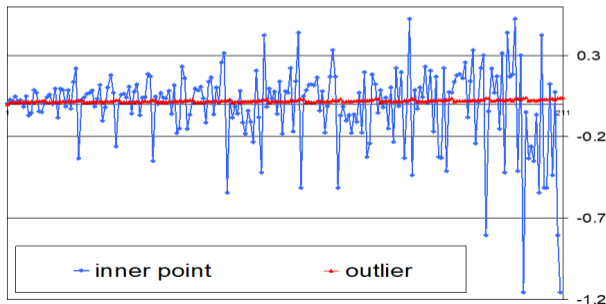
Basic Assumption

- ▶ Outliers are at the border of the data distribution
- ▶ Normal points are in the center of the data distribution

Angle-Based Approach

Model

- ▶ Consider for a given point p the angle between $\vec{p_x}$ and $\vec{p_y}$ for any two x, y from the database
- ▶ Measure the variance of the angle spectrum



Angle-Based Approach

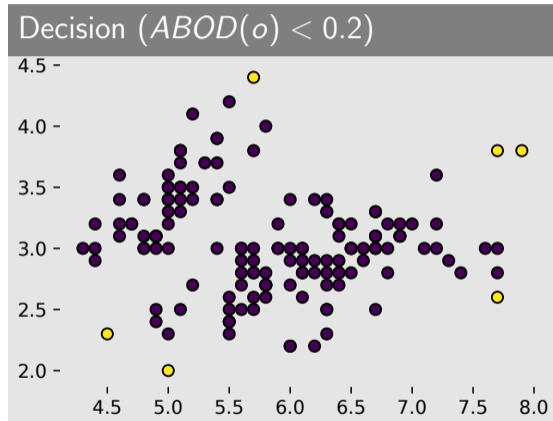
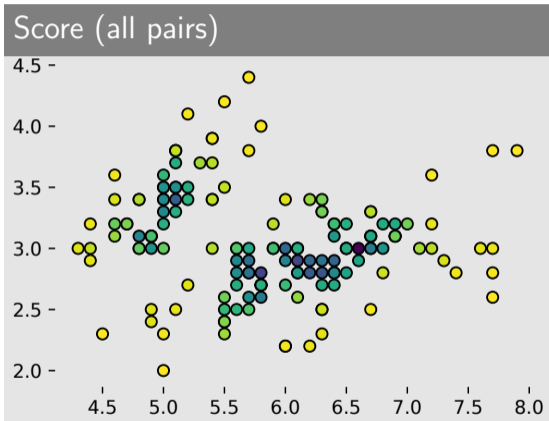
Model (cont'd)

- ▶ Weighted by the corresponding distances (for lower dimensional data sets where angles are less reliable)

$$ABOD(p) = \text{VAR}_{x,y \in D} \left(\frac{1}{\|\vec{x}_p\|_2 \|\vec{y}_p\|_2} \cos(\vec{x}_p, \vec{y}_p) \right) = \text{VAR}_{x,y \in D} \left(\frac{\langle \vec{x}_p, \vec{y}_p \rangle}{\|\vec{x}_p\|_2^2 \|\vec{y}_p\|_2^2} \right)$$

- ▶ Small ABOD \iff outlier

Angle-Based Approaches



Agenda

1. Introduction
2. Basics
3. Unsupervised Methods
 - 3.1 Frequent Pattern Mining
 - 3.2 Clustering
 - 3.3 **Outlier Detection**
 - 3.3.1 Clustering-based Outliers
 - 3.3.2 Statistical Outliers
 - 3.3.3 Distance-based Outliers
 - 3.3.4 Density-based Outliers
 - 3.3.5 Angle-based Outliers
 - 3.3.6 **Summary**
4. Supervised Methods
5. Advanced Topics

Summary

- ▶ Properties: global vs. local, labeling vs. scoring
- ▶ *Clustering-Based* Outliers: Identification as non-(cluster-members)
- ▶ *Statistical* Outliers: Assume probability distribution; outliers = unlikely to be generated by distribution
- ▶ *Distance-Based* Outliers: Distance to neighbors as outlier metric
- ▶ *Density-Based* Outliers: Relative density around the point as outlier metric
- ▶ *Angle-Based* Outliers: Angles between outliers and random point pairs vary only slightly