

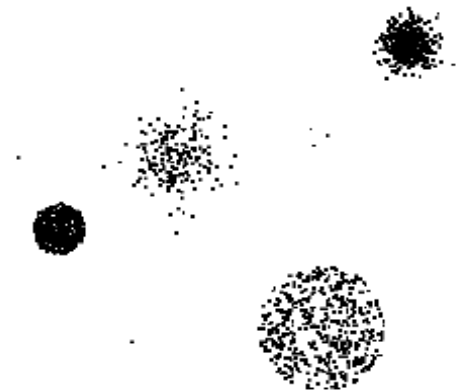
Intelligent System and Machine Learning

Outlier Detection

Slides are by Tan, Steinbach, Karpatne, Kumar

Anomaly/Outlier Detection

- What are anomalies/outliers?
 - The set of data points that are considerably different than the remainder of the data
- Natural implication is that anomalies are relatively rare
 - One in a thousand occurs often if you have lots of data
 - Context is important, e.g., freezing temps in July
- Can be important or a nuisance
 - 10 foot tall 2 year old
 - Unusually high blood pressure



Causes of Anomalies

- Data from different classes
 - Measuring the weights of oranges, but a few grapefruit are mixed in
- Natural variation
 - Unusually tall people
- Data errors
 - 200 pound 2 year old

Distinction Between Noise and Anomalies

- Noise is erroneous, perhaps random, values or contaminating objects
 - Weight recorded incorrectly
 - Grapefruit mixed in with the oranges
- Noise doesn't necessarily produce unusual values or objects
- Noise is not interesting
- Anomalies may be interesting if they are not a result of noise
- Noise and anomalies are related but distinct concepts

General Issues: Number of Attributes

- Many anomalies are defined in terms of a single attribute
 - Height
 - Shape
 - Color
- Can be hard to find an anomaly using all attributes
 - Noisy or irrelevant attributes
 - Object is only anomalous with respect to some attributes
- However, an object may not be anomalous in one attribute

General Issues: Anomaly Scoring

- ❑ Many anomaly detection techniques provide only a binary categorization
 - An object is an anomaly or it isn't
 - This is especially true of classification-based approaches
- ❑ Other approaches assign a score to all points
 - This score measures the degree to which an object is an anomaly
 - This allows objects to be ranked
- ❑ In the end, you often need a binary decision
 - Should this credit card transaction be flagged?
 - Still useful to have a score
- ❑ How many anomalies are there?

Variants of Anomaly Detection Problems

- Given a data set D , find all data points $\mathbf{x} \in D$ with anomaly scores greater than some threshold t
- Given a data set D , find all data points $\mathbf{x} \in D$ having the top- n largest anomaly scores
- Given a data set D , containing mostly normal (but unlabeled) data points, and a test point \mathbf{x} , compute the anomaly score of \mathbf{x} with respect to D

Additional Anomaly Detection Techniques

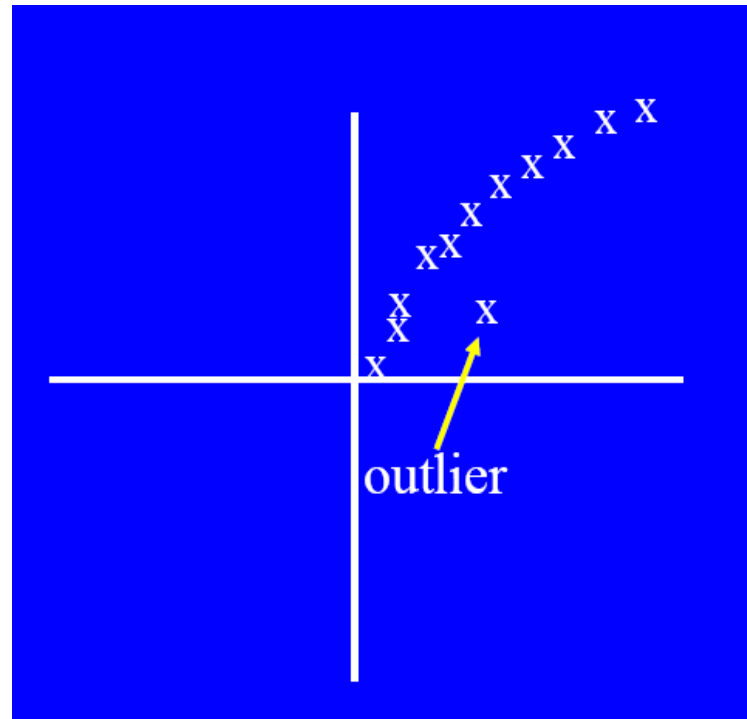
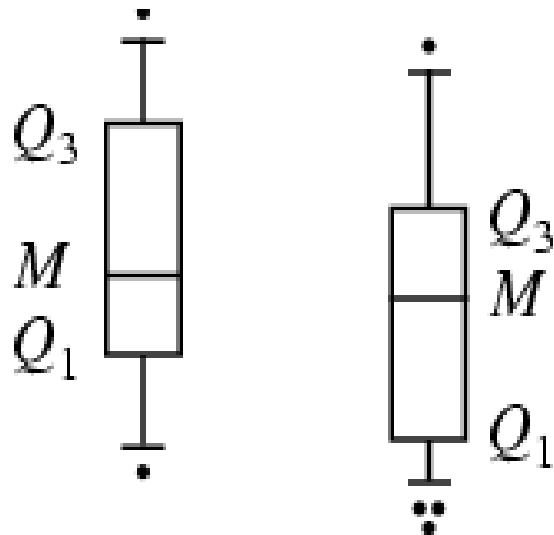
- Distance-based
 - Anomalies are points far away from other points
 - Can detect this graphically in some cases
- Density-based
 - Low density points are outliers
- Pattern matching
 - Create profiles or templates of atypical but important events or objects
 - Algorithms to detect these patterns are usually simple and efficient

Visual Approaches

□ Boxplots or scatter plots

□ Limitations

- Not automatic
- Subjective

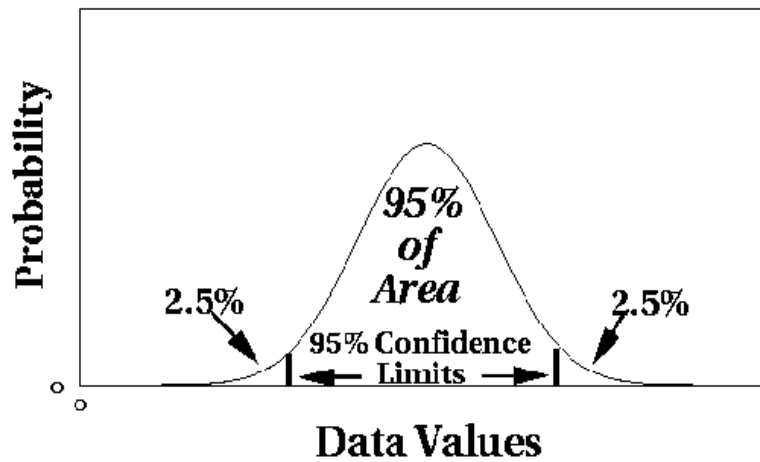


Statistical Approaches

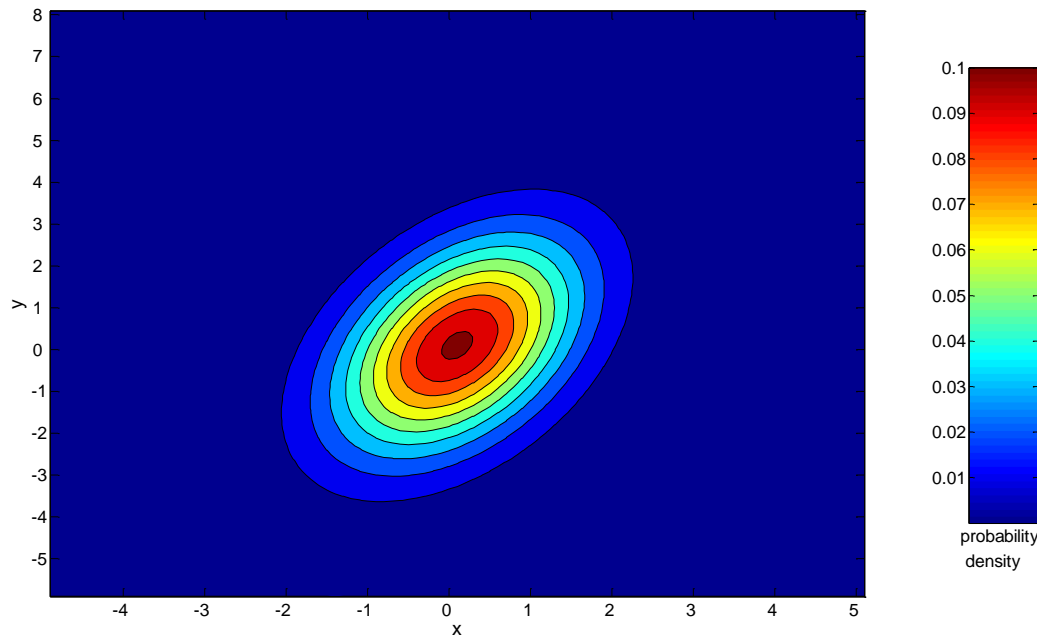
Probabilistic definition of an outlier: An outlier is an object that has a low probability with respect to a probability distribution model of the data.

- Usually assume a parametric model describing the distribution of the data (e.g., normal distribution)
- Apply a statistical test that depends on
 - Data distribution
 - Parameters of distribution (e.g., mean, variance)
 - Number of expected outliers (confidence limit)
- Issues
 - Identifying the distribution of a data set
 - ◆ Heavy tailed distribution
 - Number of attributes
 - Is the data a mixture of distributions?

Normal Distributions



**One-dimensional
Gaussian**



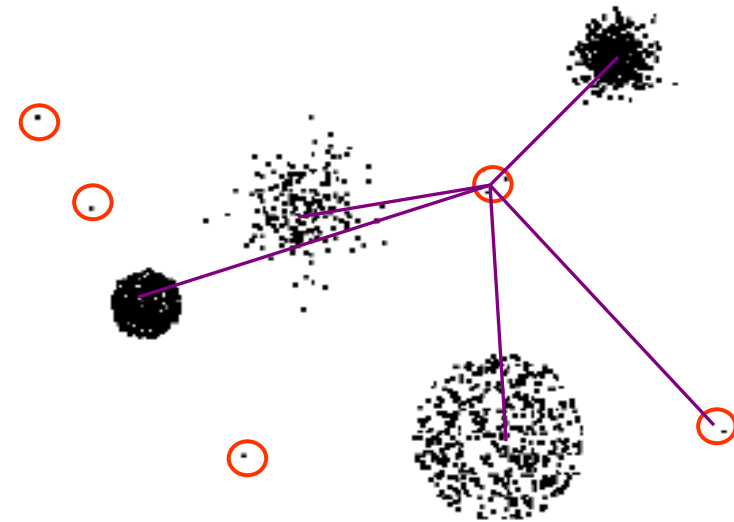
**Two-dimensional
Gaussian**

Strengths/Weaknesses of Statistical Approaches

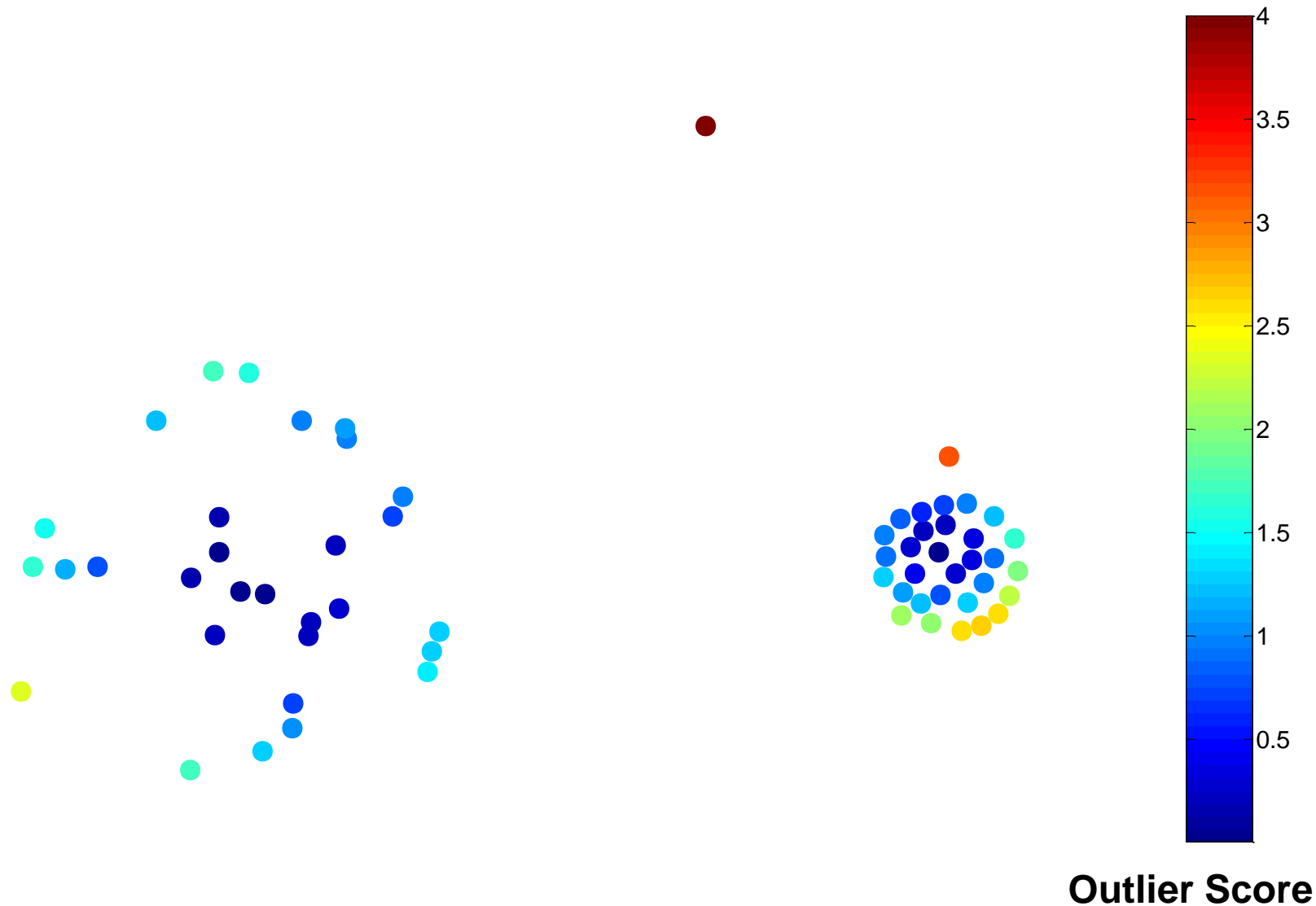
- ❑ Firm mathematical foundation
- ❑ Can be very efficient
- ❑ Good results if distribution is known
- ❑ In many cases, data distribution may not be known
- ❑ For high dimensional data, it may be difficult to estimate the true distribution
- ❑ Anomalies can distort the parameters of the distribution

Clustering-Based Approaches

- **Clustering-based Outlier:** An object is a cluster-based outlier if it does not strongly belong to any cluster
 - For prototype-based clusters, an object is an outlier if it is not close enough to a cluster center
 - For density-based clusters, an object is an outlier if its density is too low
 - For graph-based clusters, an object is an outlier if it is not well connected
- Other issues include the impact of outliers on the clusters and the number of clusters



Relative Distance of Points from Closest Centroid



Strengths/Weaknesses of Distance-Based Approaches

- Simple method
- Many clustering techniques can be used
- Can be difficult to decide on a clustering technique, can be difficult to decide on number of clusters
- Outliers can distort the clusters