



# 数据挖掘导论

## Introduction to Data Mining

### Outlier Detection



数据智能实验室  
DATA INTELLIGENCE LABORATORY



浙江大学  
Zhejiang University

# Agenda

## □ Concepts and Applications

## □ Statistical Methods

## □ Graphical Method

## □ Density-Based Methods

## □ Isolation Tree

## □ Summary

# Outlier Detection

## □ Also called anomaly detection

- ✓ Identify objects that are different from most other objects

## □ Causes of anomalies

- ✓ Data from different class of object or underlying mechanism, e.g., rare disease, fraud detection
- ✓ Natural variation e.g., tails on a Gaussian distribution
- ✓ Data measurement or collection errors

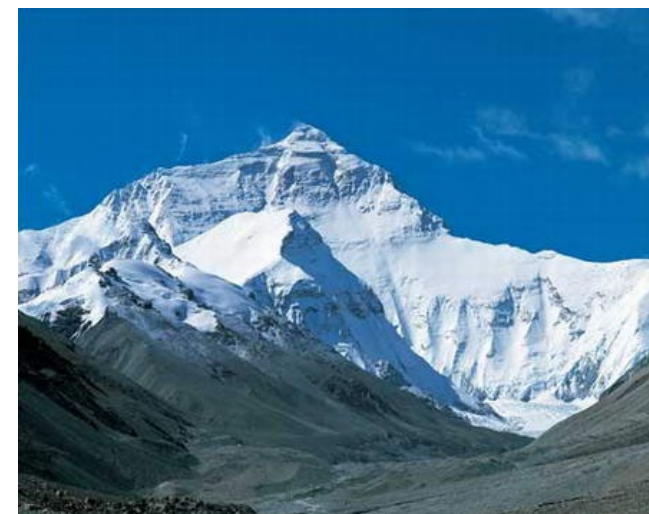
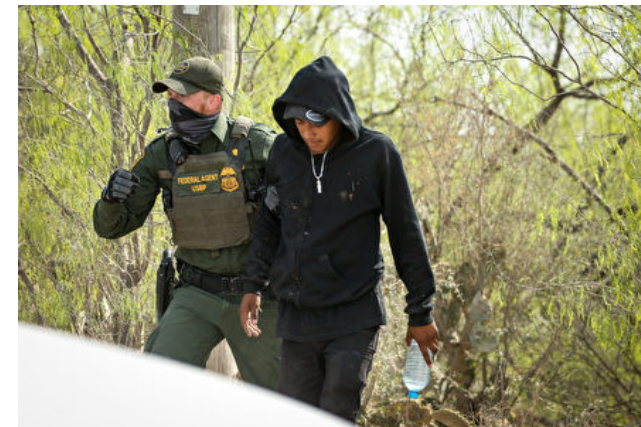




# Applications – Spam/Fraud Detection



# Applications – Surveillance Outlier Event

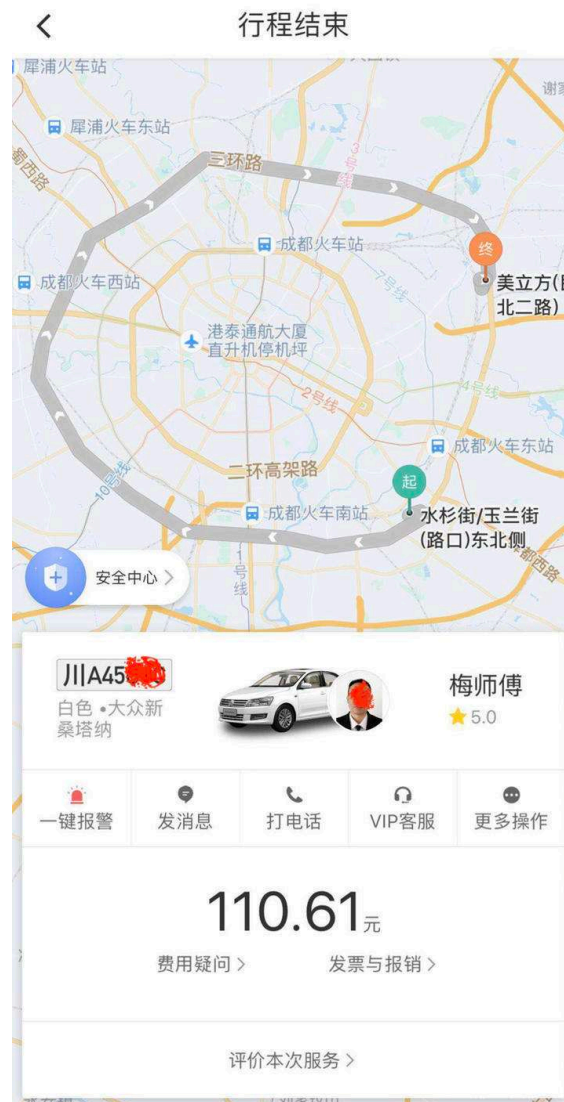




# Applications – Detour Detection



一个多小时行程结束以后，吴小姐傻眼了，滴滴系统发来的账单显示吴小姐要支付**540多万元**的打车费用，而打车时间显示的则是**1000多万分钟**。



图片来源

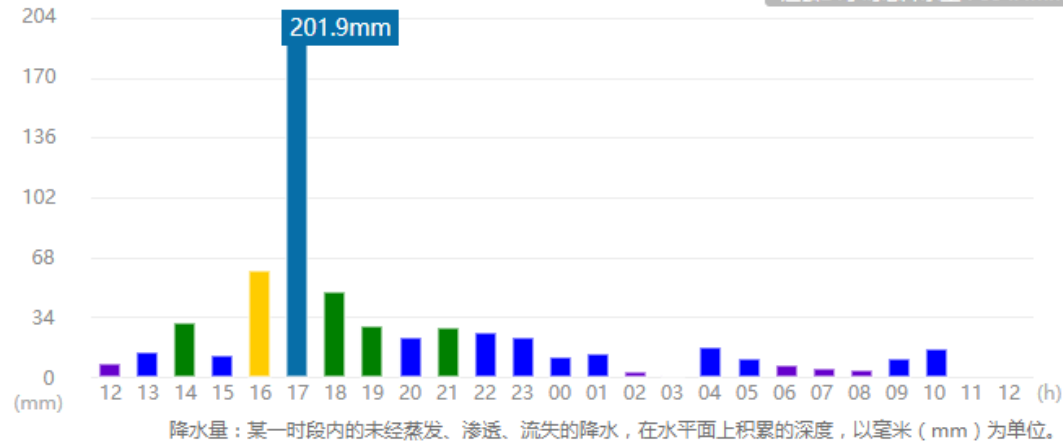
醉酒乘客打滴滴回家 司机疯狂绕路带其"环游"1小时

# Applications – Time Series Outlier Detection

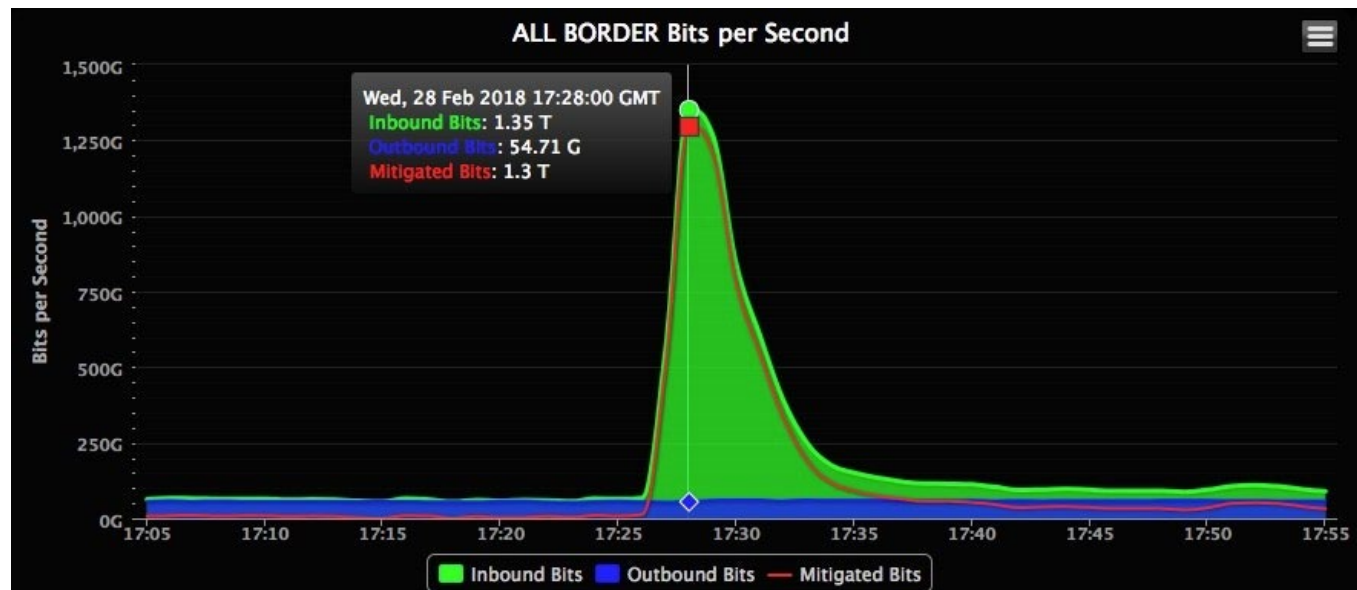
整点天气实况

空气质量 | 温度 | 相对湿度 | 降水量 | 风力风向

过去24小时总降水量：594.4mm



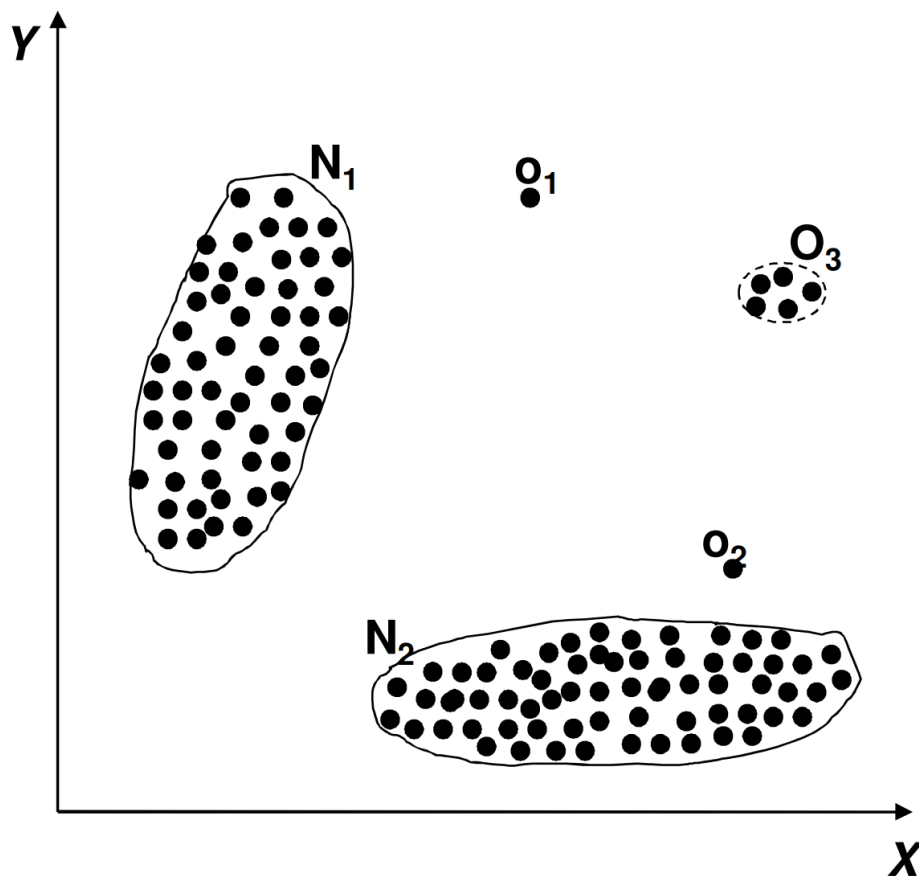
国家防总：7月20日郑州最大小时降雨量达201.9毫米突破历史极值



GitHub Survived the Biggest DDoS Attack Ever Recorded

# Point Anomaly

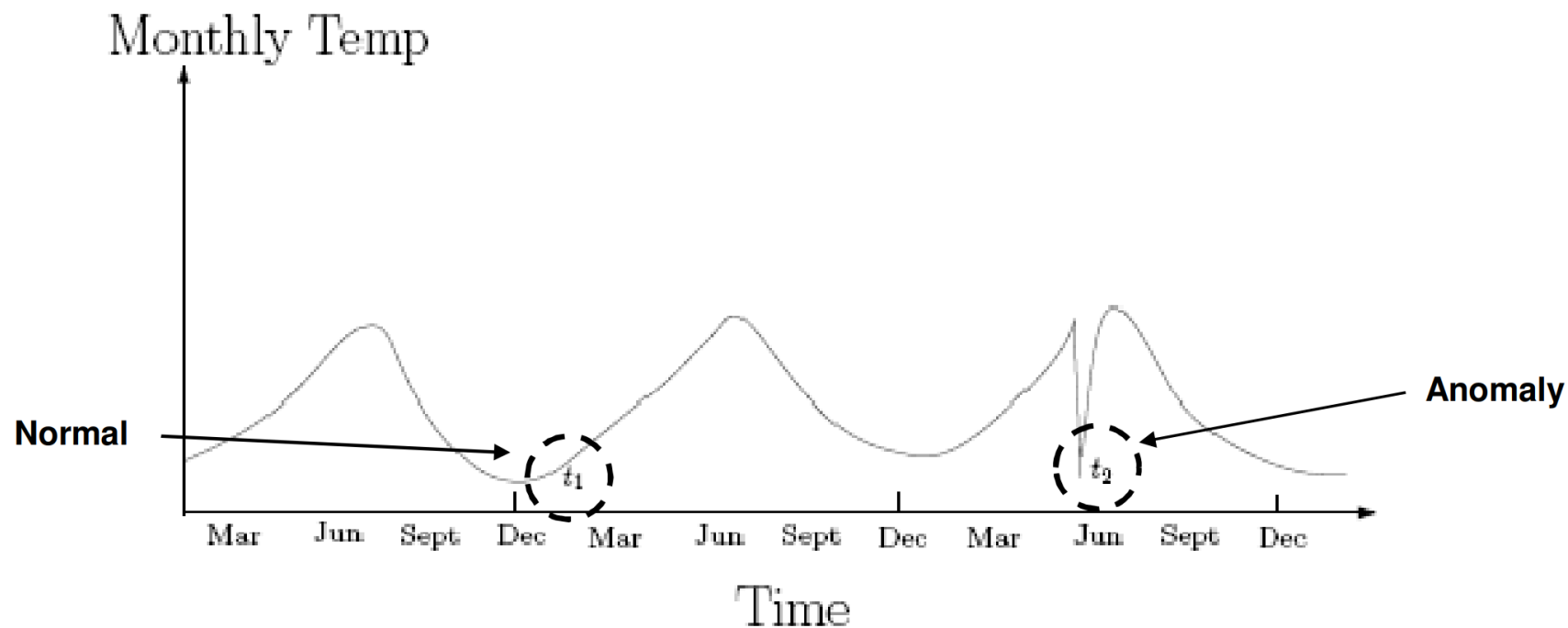
- An individual data instance is anomalous with respect to the data





# Contextual Anomaly

- An individual data instance is anomalous within a context
- Dinosaurs are common in cretaceous, but anomalous nowadays



# Agenda

- Concepts and Applications

- **Statistical Methods**

- Graphical Method

- Density-Based Methods

- Isolation Tree

- Summary

# Statistical Outlier Detection

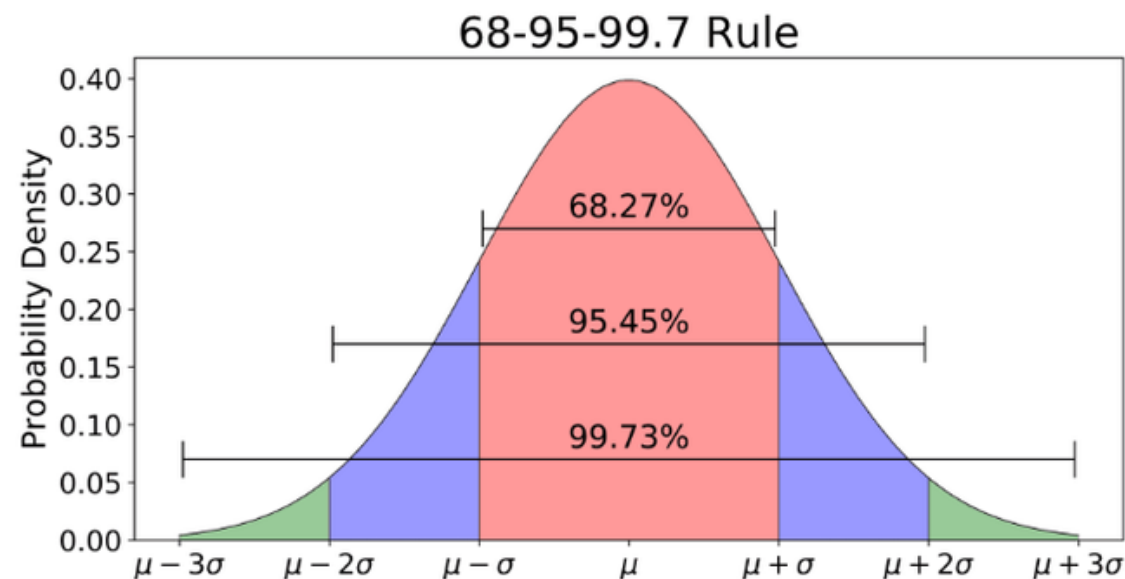
## □ Core idea

- ✓ Fit the data with a probabilistic model
- ✓ Outliers are instances with low probability

## □ Univariate Gaussian Distribution

- ✓ E.g., the height of students in ZJU
- ✓ Outlier defined by z-score > threshold

$$z = \frac{x - \mu}{\sigma}$$



The 68-95-99.7 Rule for a Normal Distribution

# Grubbs' Test

❑ **Example: Given sampled data points 5, 10, 9.5, 9.8, 9.9, can we reject the outlier with 95% confidence?**

- ✓ Find the mean ( $\bar{x}$ ) and standard deviation of the data set ( $\bar{x}=8.84$ ,  $\sigma=2.119$ )
- ✓ Calculate z-score =  $|5-8.84|/2.119=1.812$
- ✓ Calculate the G-critical value, usually lookup from table (1.67 in this example)
- ✓ Reject the outlier if the test statistic is greater than G-critical value

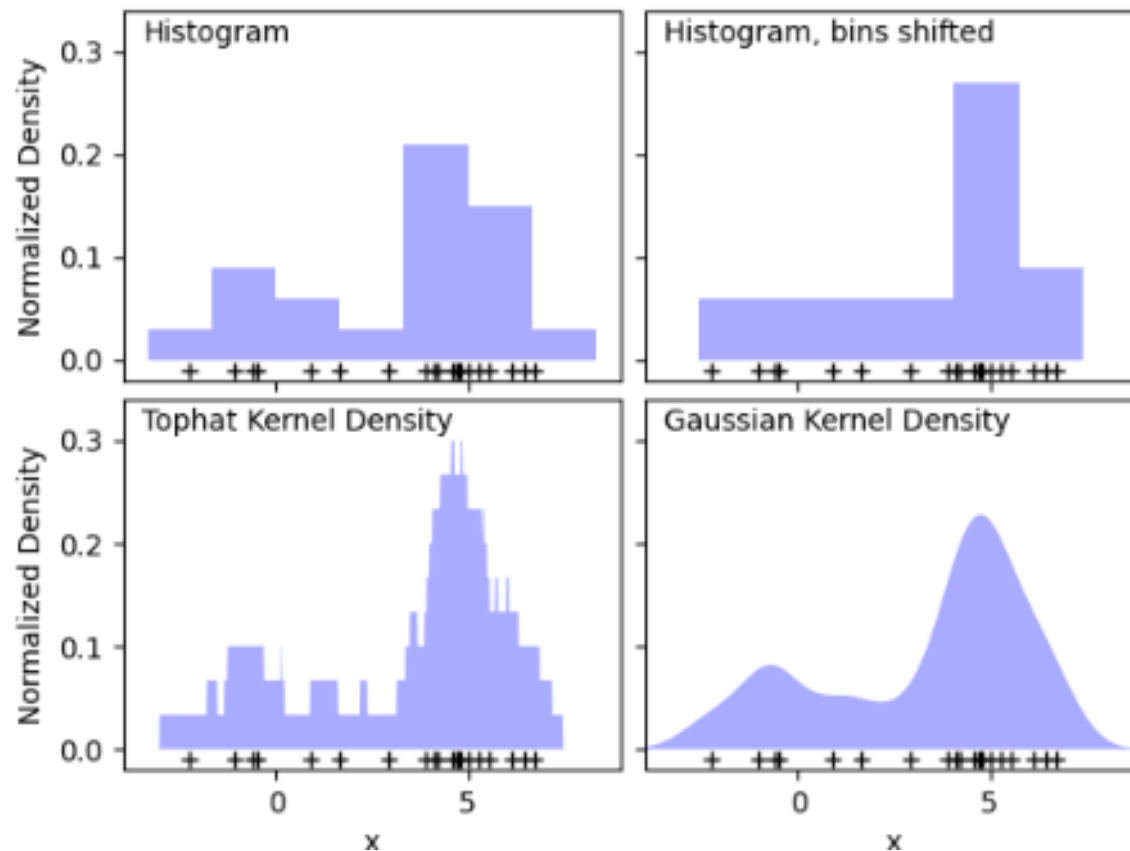
Alpha					
N	0.1	0.075	0.05	0.025	0.01
3	1.15	1.15	1.15	1.15	1.15
4	1.42	1.44	1.46	1.48	1.49
5	1.6	1.64	1.67	1.71	1.75
6	1.73	1.77	1.82	1.89	1.94
7	1.83	1.88	1.94	2.02	2.1
8	1.91	1.96	2.03	2.13	2.22
9	1.98	2.04	2.11	2.21	2.32



# Kernel Density Estimation

□ Estimate the PDF of a random variable in a non-parametric way

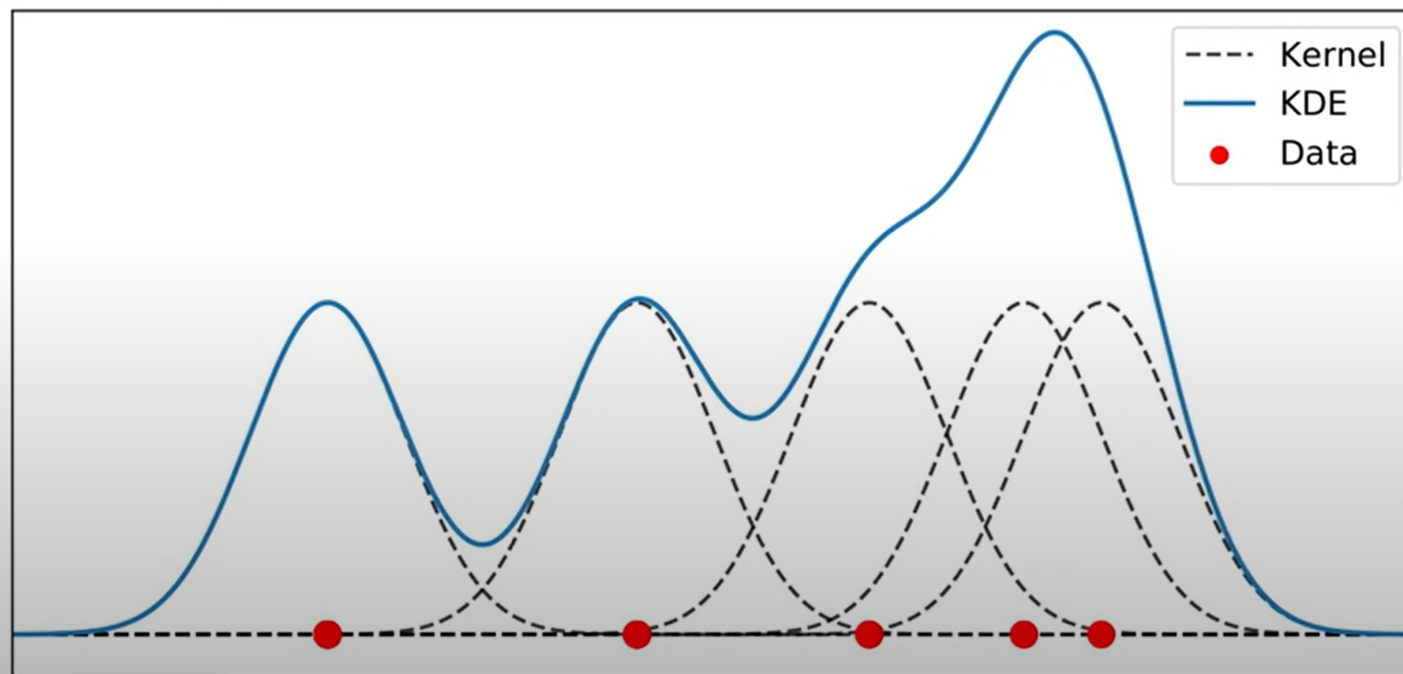
✓ It's related to a histogram but with a data smoothing technique.



# Kernel Density Estimation

On every data point  $x_i$ , we place a kernel function  $K$ . The kernel density estimate is

$$\hat{f}(x) = \frac{1}{N} \sum_{i=1}^N K(x - x_i).$$



# Agenda

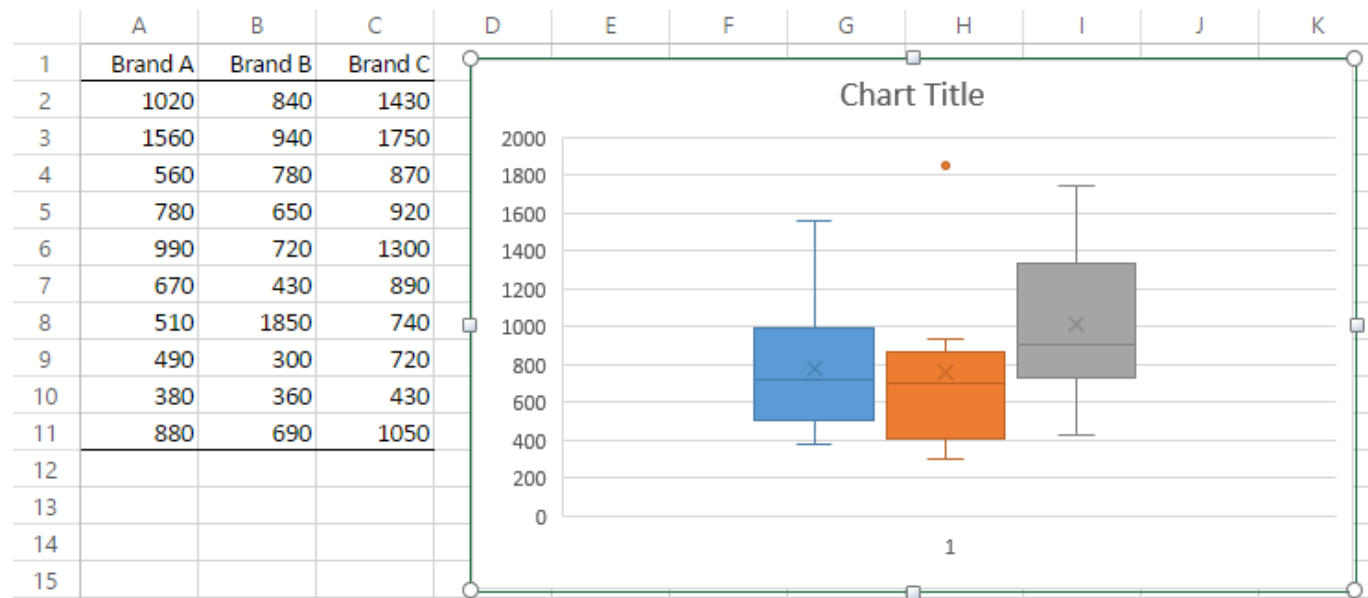
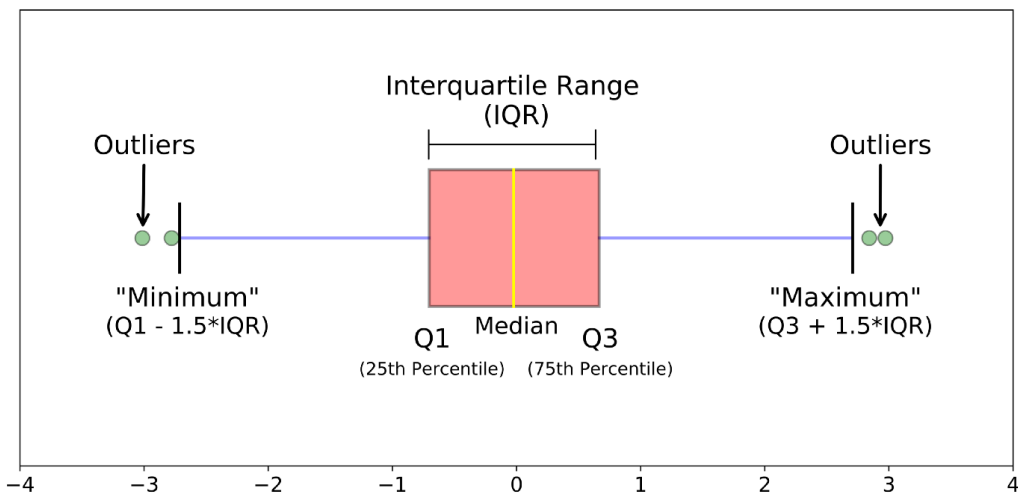
- Concepts and Applications
- Statistical Methods
- Graphical Method
- Density-Based Methods
- Isolation Tree
- Summary

# Graphical Outlier Detection

## □ Graphical approach to outlier detection

- ✓ Look at a plot of the data
- ✓ Human decides if data is an outlier

## □ Box plot

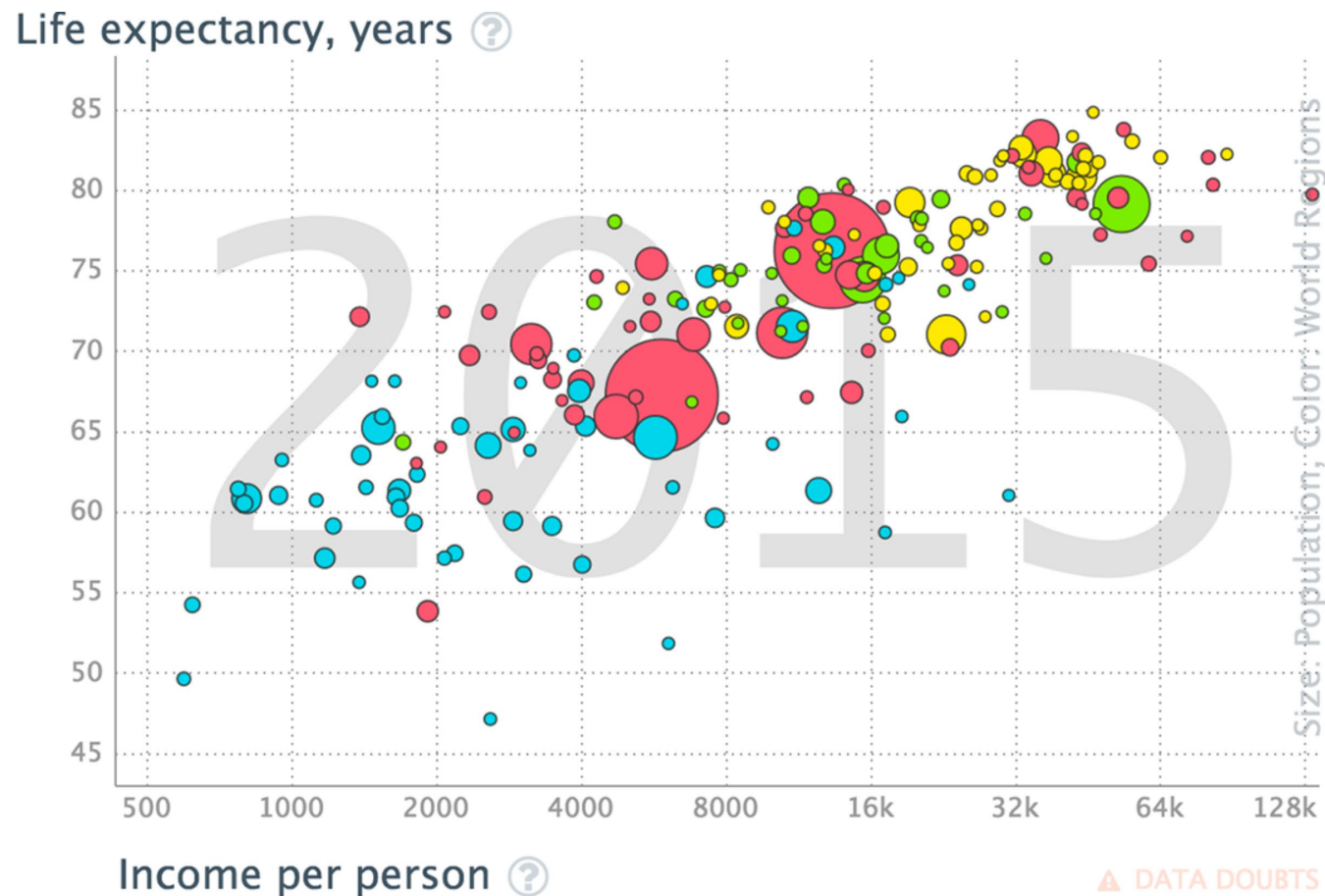




# Graphical Outlier Detection

## □ Scatter Plot Array

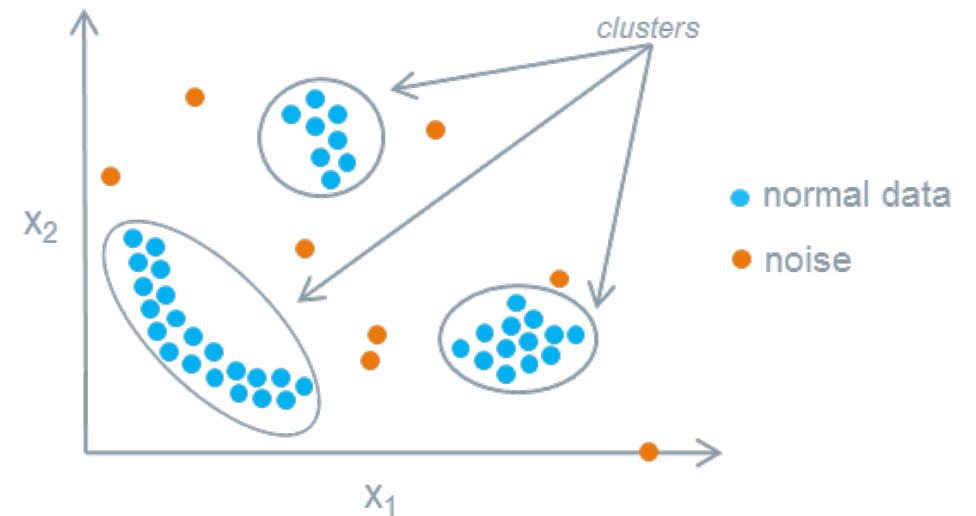
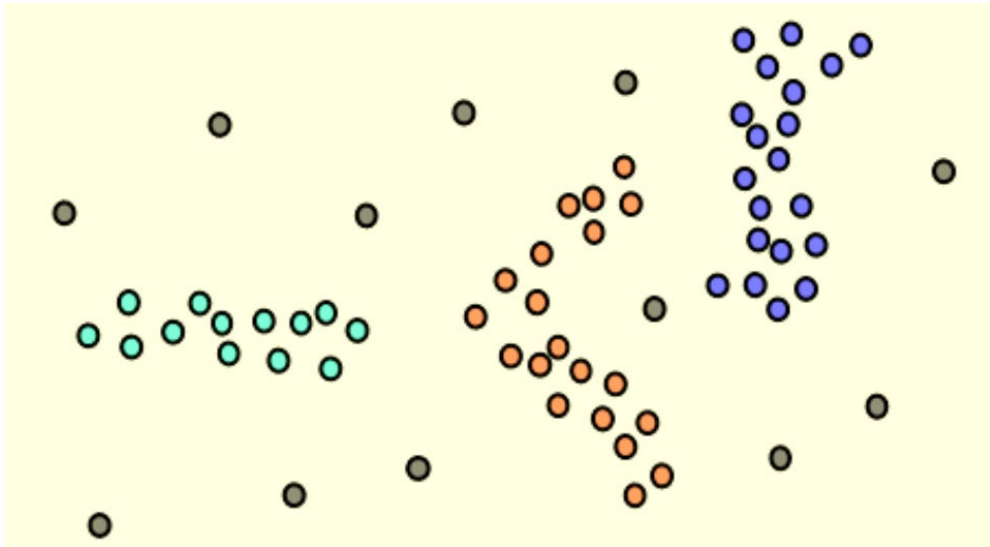
- ✓ Shows how multiple variables are related. The matrix can also identify outliers in multiple scatter plots



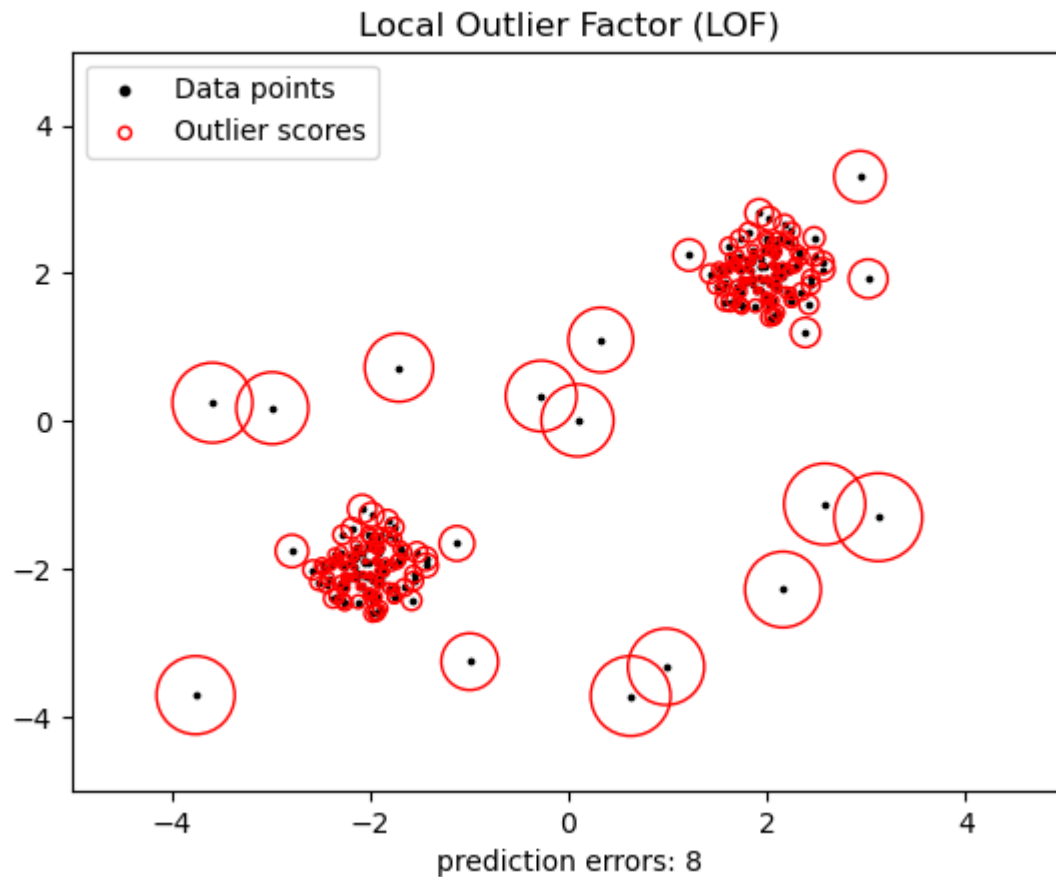
# Agenda

- Concepts and Applications
- Statistical Methods
- Graphical Method
- **Density-Based Methods**
- Isolation Tree
- Summary

# DBSCAN



# Local Outlier Factor (LOF)





# Local Outlier Factor (LOF)

## □ K-distance(x)

- ✓ The distance of object x to its k-th nearest neighbor

## □ Reachability Distance (RD)

- ✓ The maximum of the distance of two points and the k-distance of the second point
- ✓  $RD(a, b) = \max\{k\text{-distance}(b), \text{distance}(a, b)\}$

## □ Local Reachability Density (LRD)

- ✓ It refers to how far we need to go from the point we are at to reach the next point or set of points.
- ✓  $LRD(a) = 1 / (\sum(RD(a, n)) / k)$ , where n is a member in the k nearest neighbors

# Local Outlier Factor (LOF)

## □ Local Outlier Factor

- ✓ Average ratio of LRD of neighbors of  $p$  and LRD of  $p$
- ✓  $LOF(p) = [ (LRD(1st. neighbor) + LRD(2nd. neighbor) + \dots + LRD(kth. neighbor)) / LRD(p) ] / k$

## □ Property of LOF

- ✓ A point is an outlier if its  $LOF \gg 1$

$LOF(k) \sim 1$  means **Similar density as neighbors**,

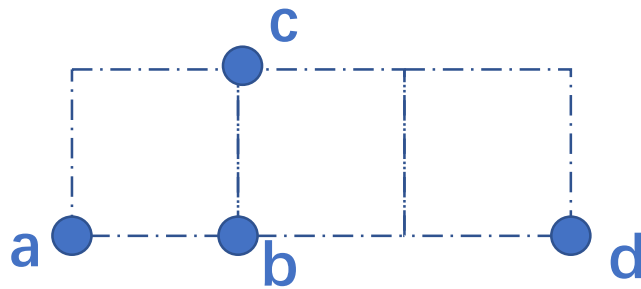
$LOF(k) < 1$  means **Higher density than neighbors**

$LOF(k) > 1$  means **Lower density than neighbors**

# Local Outlier Factor (LOF)

## □ Example of LOF calculation

- ✓ Given 4 points  $a(0,0)$ ,  $b(0,1)$ ,  $c(1,1)$  and  $d(3,0)$
- ✓ Let  $k=2$  and use Manhattan distance
- ✓ Calculate the LOF for each point and show the top-1 outlier



# Local Outlier Factor (LOF)

□ Step1: calculate all the distances between each two data points

$$\text{dist}(a, b) = 1$$

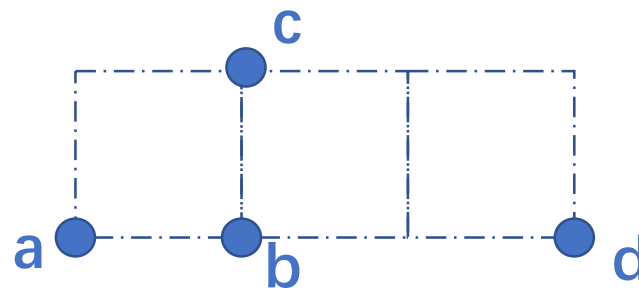
$$\text{dist}(a, c) = 2$$

$$\text{dist}(a, d) = 3$$

$$\text{dist}(b, c) = 1$$

$$\text{dist}(b, d) = 3+1=4$$

$$\text{dist}(c, d) = 2+1=3$$





# Local Outlier Factor (LOF)

□ Step2: calculate k-distance for each object

$$N_2(a) = \{b, c\}$$

$$N_2(b) = \{a, c\}$$

$$N_2(c) = \{b, a\}$$

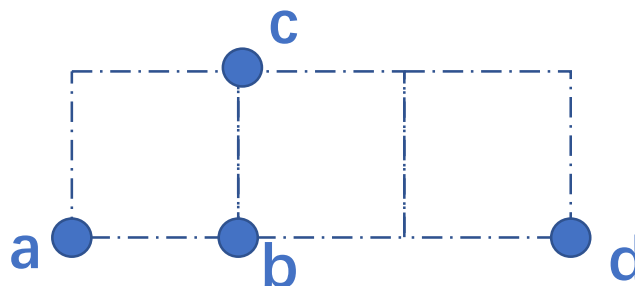
$$N_2(d) = \{a, c\}$$

$$\text{dist}_2(a) = \text{dist}(a, c) = 2 \text{ (c is the 2}^{nd} \text{ nearest neighbor)}$$

$$\text{dist}_2(b) = \text{dist}(b, a) = 1 \text{ (a/c is the 2}^{nd} \text{ nearest neighbor)}$$

$$\text{dist}_2(c) = \text{dist}(c, a) = 2 \text{ (a is the 2}^{nd} \text{ nearest neighbor)}$$

$$\text{dist}_2(d) = \text{dist}(d, a) = 3 \text{ (a/c is the 2}^{nd} \text{ nearest neighbor)}$$



# Local Outlier Factor (LOF)

□ Step3: calculate RD and LRD for all the objects

$$RD(a,b)=1$$

$$RD(a,c)=2$$

$$LRD(a)=2/(1+2)=0.667$$

$$RD(c,a)=2$$

$$RD(c,b)=1$$

$$LRD(c)=2/(1+2)=0.667$$

$$RD(b,a)=2$$

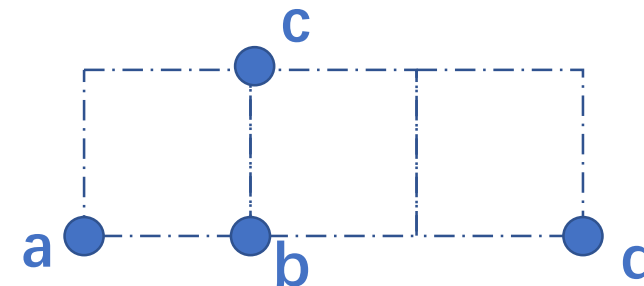
$$RD(b,c)=2$$

$$LRD(b)=2/(2+2)=0.5$$

$$RD(d,a)=3$$

$$RD(d,c)=3$$

$$LRD(d)=2/(3+3)=0.333$$



# Local Outlier Factor (LOF)

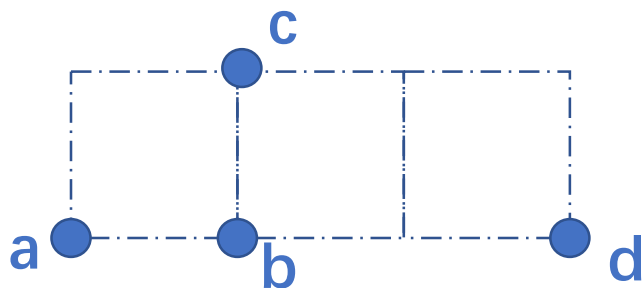
□ Step4: calculate LOF for all the objects

$$\text{LOF}(a) = [(0.5 + 0.667) / 0.667] / 2 = 0.875$$

$$\text{LOF}(b) = [(0.667 + 0.667) / 0.5] / 2 = 1.334$$

$$\text{LOF}(c) = [(0.667 + 0.5) / 0.667] / 2 = 0.875$$

$$\text{LOF}(d) = [(0.667 + 0.667) / 0.333] / 2 = 2 \longrightarrow \text{outlier}$$



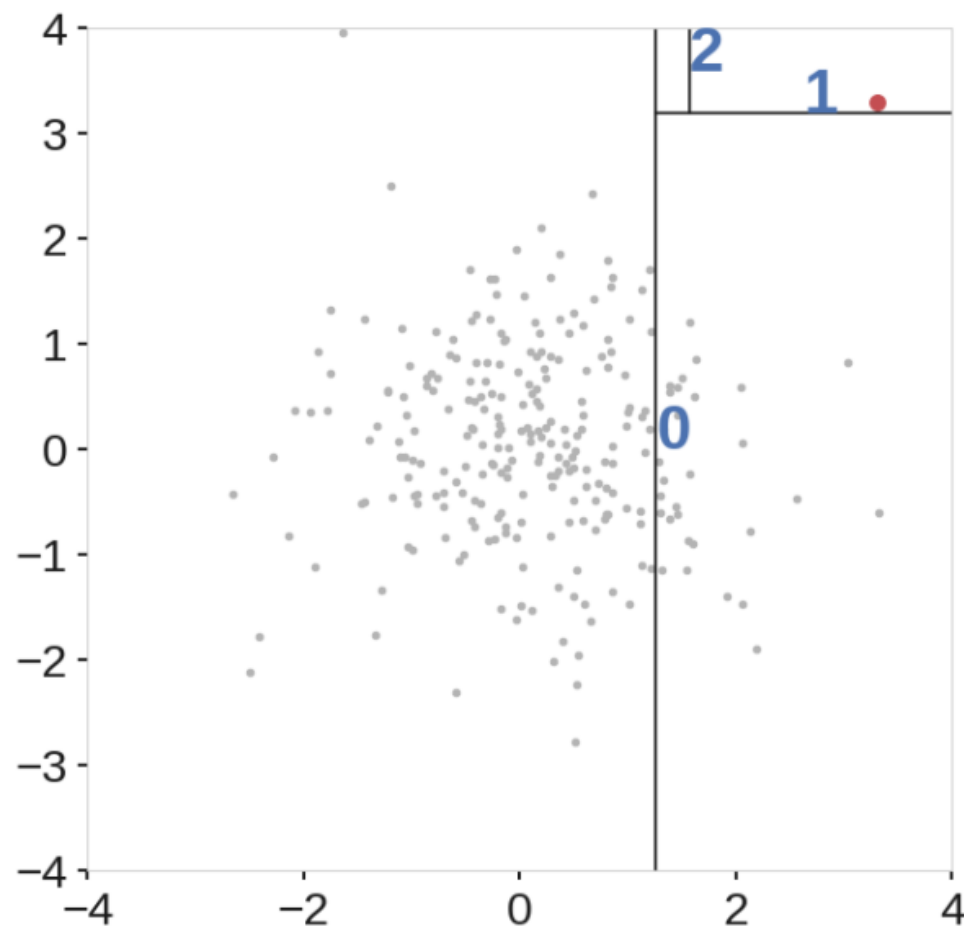
# Agenda

- Concepts and Applications
- Statistical Methods
- Graphical Method
- Density-Based Methods
- Isolation Tree**
- Summary

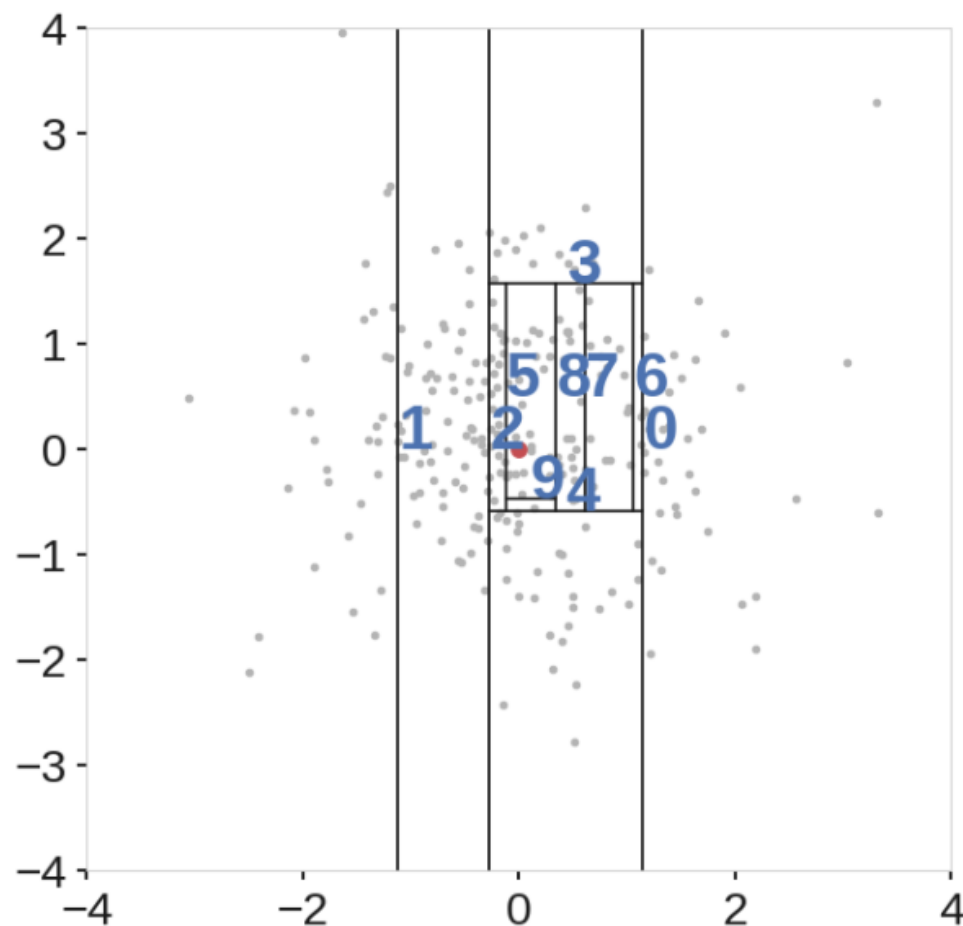
# Isolation Tree

- ❑ **An unsupervised learning algorithm based on decision tree**
  - ✓ Builds decision trees with random partitioning
- ❑ **Build a decision tree with random partitioning**
  - ✓ Randomly select a feature from the given set of features
  - ✓ Randomly select a split value between the max and min values of that feature
- ❑ **Such random partitioning produces shorter paths in trees for the anomalous data points**
- ❑ **When a forest of random trees collectively produce shorter path lengths for particular samples, they are highly likely to be anomalies**

# Isolation Tree



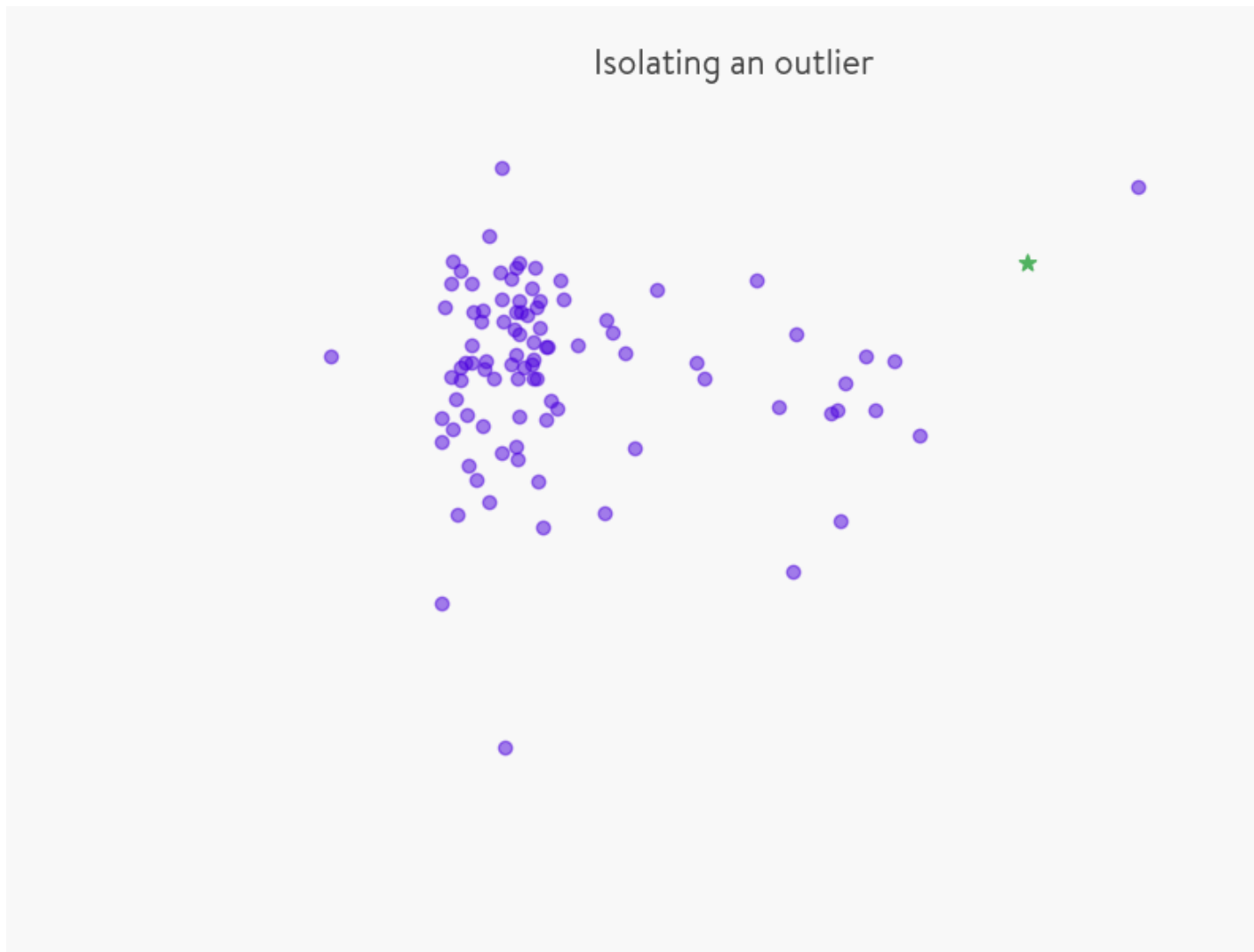
(a) Anomaly point



(b) Nominal point



# Isolation Tree



# Agenda

- Concepts and Applications
- Statistical Methods
- Graphical Method
- Density-Based Methods
- Isolation Tree
- Summary