

CPSC 340: Machine Learning and Data Mining

Outlier Detection

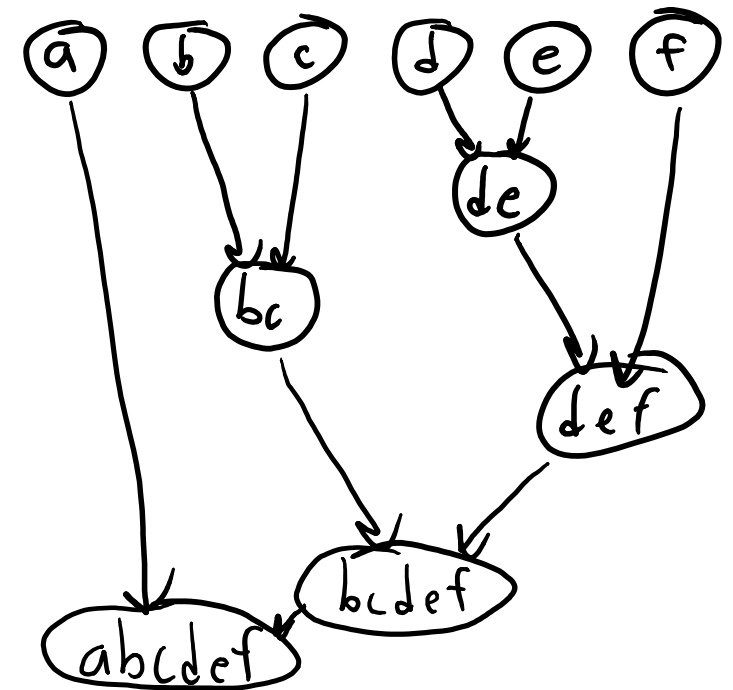
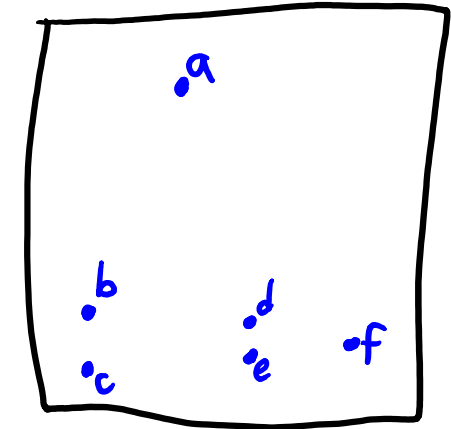
Fall 2018

Admin

- **Assignment 2** is due Friday.
- Assignment 1 grades available?
- Midterm rooms are now booked.
 - October 18th at 6:30pm (BUCH A102 and A104).
- Mike and I will get a little out of sync over the next few lectures.
 - Keep this in mind if you alternating between our lectures.

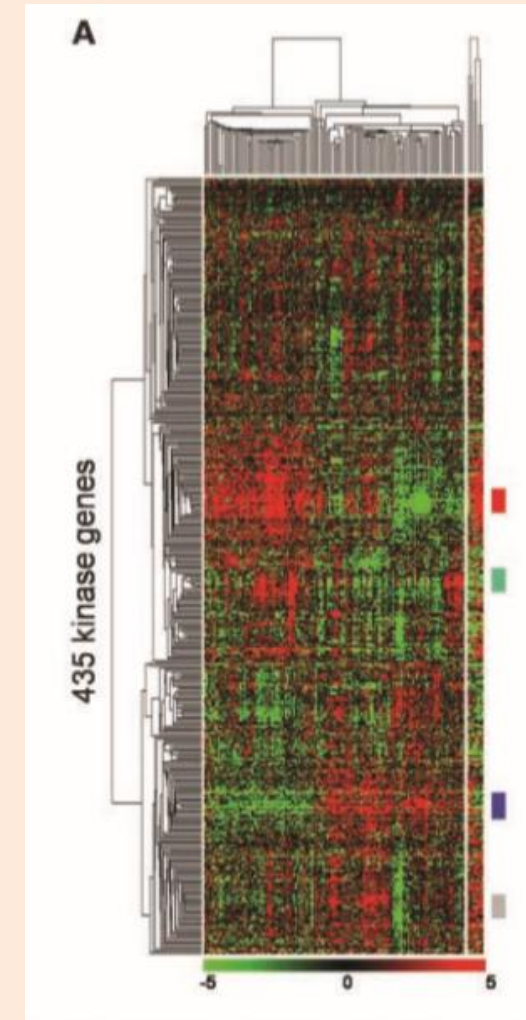
Last Time: Hierarchical Clustering

- We discussed **hierarchical clustering**:
 - Perform **clustering at multiple scales**.
 - Output is usually a **tree diagram** (“dendrogram”).
 - Reveals much more structure in data.
 - Usually non-parametric:
 - At finest scale, every point is its own clusters.
- We discussed some application areas:
 - Animals (phylogenetics).
 - Languages.
 - Stories.
 - Fashion.



Application: Medical data

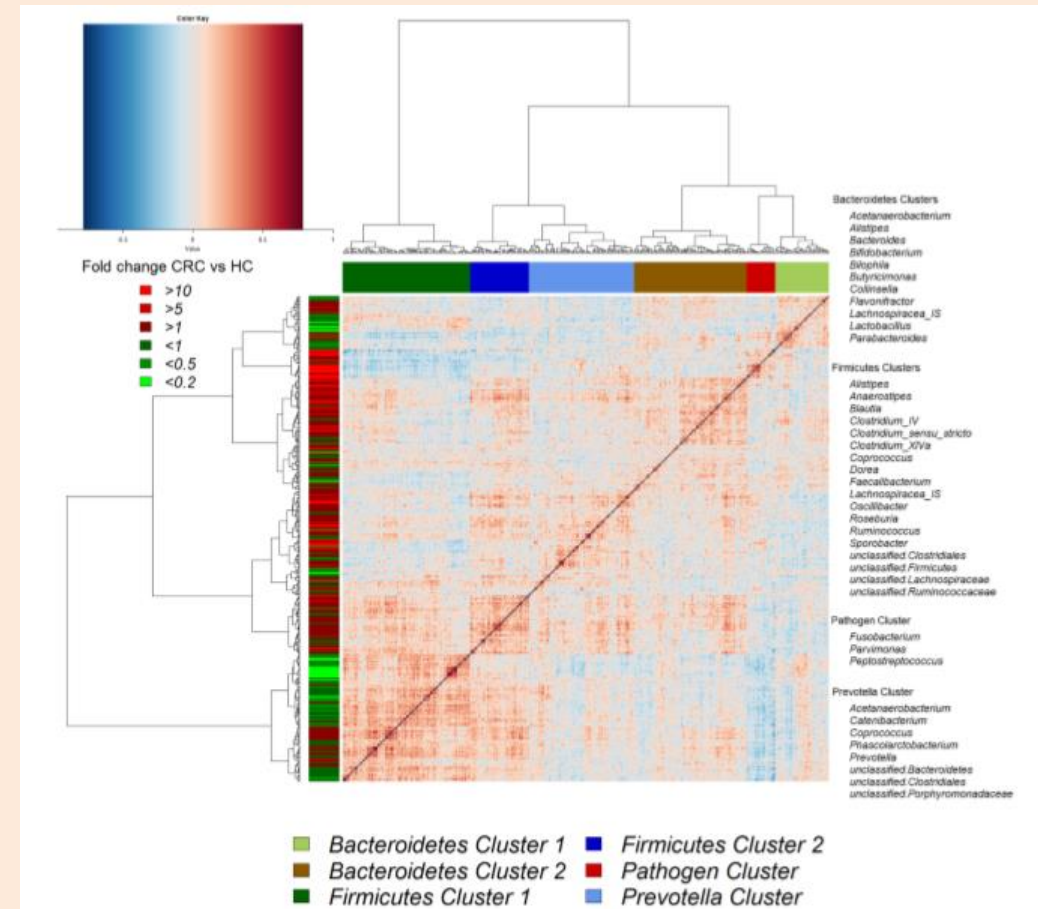
- Hierarchical clustering is very common in **medical data analysis**.
 - Clustering different samples of breast cancer:
 - Note: they are plotting X^T (samples are columns).
 - They've sorted the columns to make the plot look nicer.
 - Notice they also clustered and sorted the features (rows).
 - Gives information about **relationship between features**.



Application: Medical data

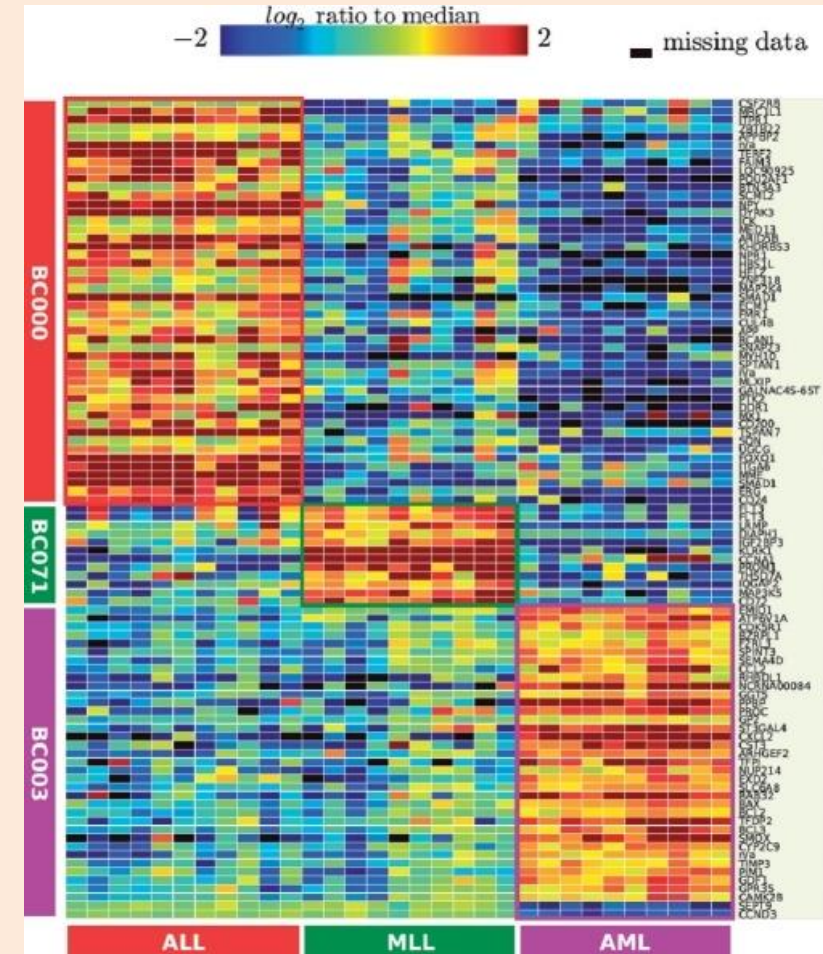
- Hierarchical clustering is very common in **medical data analysis**.
 - Clustering different samples of colorectal cancer:

- Note: the matrix is 'n' by 'n'.
 - Each matrix element gives correlation.
 - Clusters should look like “blocks” on diagonal.
 - Order of examples is reversed in columns.
 - This is why diagonal goes from bottom-to-top.



Other Clustering Methods

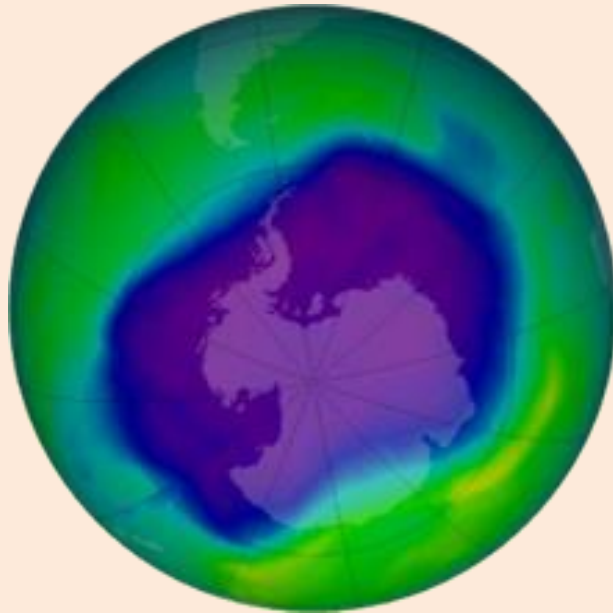
- Mixture models:
 - Probabilistic clustering.
- Mean-shift clustering:
 - Finds local “modes” in density of points.
 - Alternative approach to vector quantization.
- Bayesian clustering:
 - A variant on ensemble methods.
 - Averages over models/clustering, weighted by “prior” belief in the model/clustering.
- Biclustering:
 - Simultaneously cluster examples and features.
- Spectral clustering and graph-based clustering:
 - Clustering of data described by graphs.



(pause)

Motivating Example: Finding Holes in Ozone Layer

- The huge Antarctic ozone hole was “discovered” in 1985.

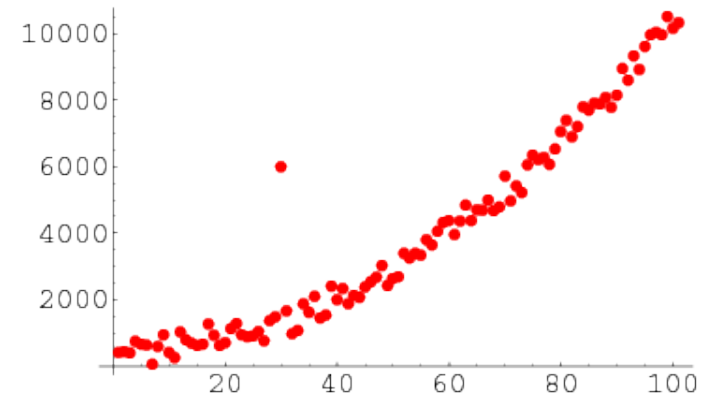
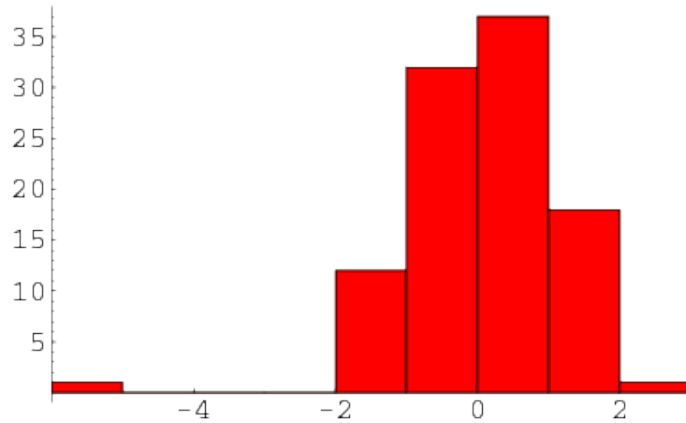


- It had been in satellite data since 1976:
 - But it was flagged and filtered out by quality-control algorithm.

Outlier Detection

- **Outlier detection:**

- Find observations that are “unusually different” from the others.
- Also known as “anomaly detection”.
- May want to remove outliers, or be interested in the outliers themselves (security).




- **Some sources of outliers:**

- Measurement errors.
- Data entry errors.
- Contamination of data from different sources.
- Rare events.

Applications of Outlier Detection

- Data cleaning.
- Security and fault detection (network intrusion, DOS attacks).
- Fraud detection (credit cards, stocks, voting irregularities).

Transaction Date	Posted Date	Transaction Details	Debit	Credit
Aug. 27, 2015	Aug. 28, 2015	 BEAN AROUND THE WORLD VANCOUVER, BC	\$10.95	

- Detecting natural disasters (underwater earthquakes).
- Astronomy (find new classes of stars/planets).
- Genetics (identifying individuals with new/ancient genes).

Classes of Methods for Outlier Detection

1. Model-based methods.
 2. Graphical approaches.
 3. Cluster-based methods.
 4. Distance-based methods.
 5. Supervised-learning methods.
- Warning: this is the topic with the most ambiguous “solutions”.

But first...

- Usually it's good to do some **basic sanity checking**...

Egg	Milk	Fish	Wheat	Shellfish	Peanuts	Peanuts	Sick?
0	0.7	0	0.3	0	0	0	1
0.3	0.7	0	0.6	-1	3	3	1
0	0	0	"sick"	0	1	1	0
0.3	0.7	1.2	0	0.10	0	0.01	2
900	0	1.2	0.3	0.10	0	0	1

- Would any values in the column cause a Python/Julia **"type" error**?
- What is the **range of numerical features**?
- What are the **unique entries for a categorical feature**?
- Does it look like parts of the table are **duplicated**?
- These types of simple errors are VERY common in real data.

Model-Based Outlier Detection

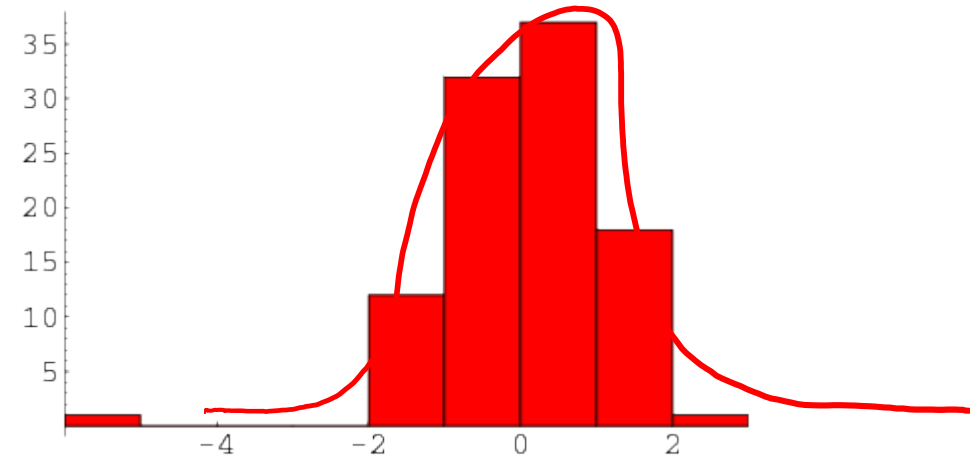
- Model-based outlier detection:
 1. Fit a probabilistic model.
 2. Outliers are examples with low probability.

- Example:

- Assume data follows normal distribution.
- The z-score for 1D data is given by:

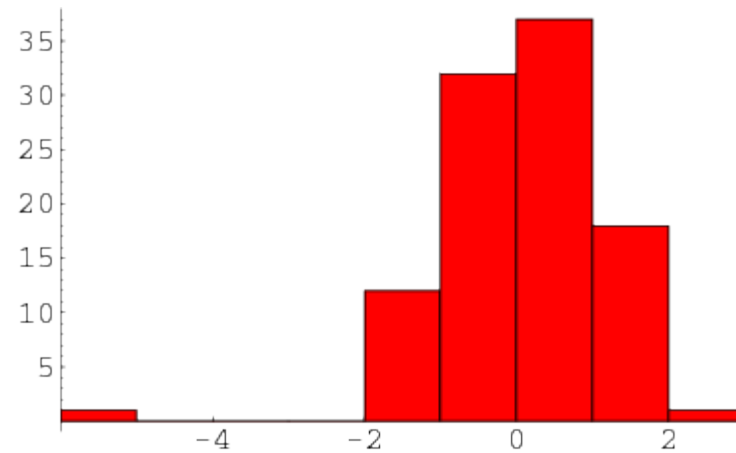
$$z_i = \frac{x_i - \mu}{\sigma} \quad \text{where } \mu = \frac{1}{n} \sum_{i=1}^n x_i \quad \text{and } \sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

- “Number of standard deviations away from the mean”.
- Say “outlier” if $|z| > 4$, or some other threshold.

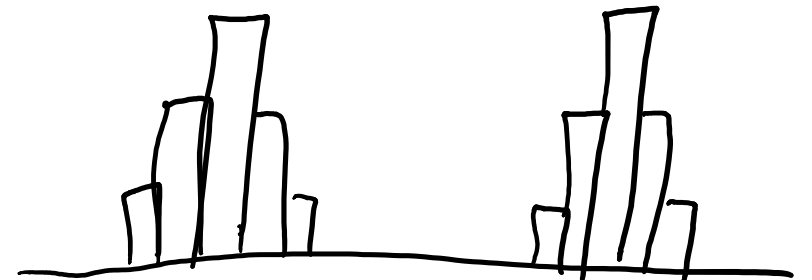


Problems with Z-Score

- Unfortunately, the **mean and variance are sensitive to outliers.**



- Possible fixes: **use quantiles, or sequentially remove worse outlier.**
- The z-score also assumes that data is “uni-modal”.
 - Data is concentrated around the mean.



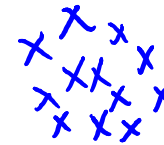
Global vs. Local Outliers

- Is the **red point** an outlier?



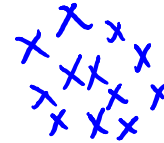
Global vs. Local Outliers

- Is the **red point** an outlier? What if we add the **blue points**?



Global vs. Local Outliers

- Is the **red point** an outlier? What if we add the **blue points**?



- Red point has the **lowest z-score**.
 - In the first case it was a “**global**” outlier.
 - In this second case it’s a “**local**” outlier:
 - Within normal data range, but **far from other points**.
- It’s hard to precisely define “outliers”.

Global vs. Local Outliers

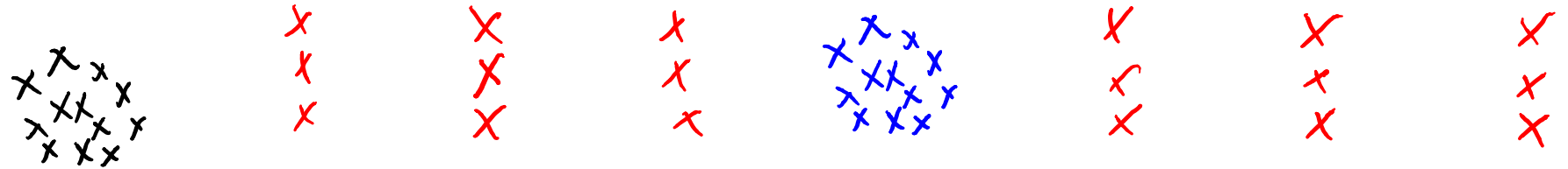
- Is the **red point** an outlier? What if we add the **blue points**?



- Red point has the **lowest z-score**.
 - In the first case it was a “**global**” outlier.
 - In this second case it’s a “**local**” outlier:
 - Within normal data range, but **far from other points**.
- It’s hard to precisely define “outliers”.
 - Can we have **outlier groups**?

Global vs. Local Outliers

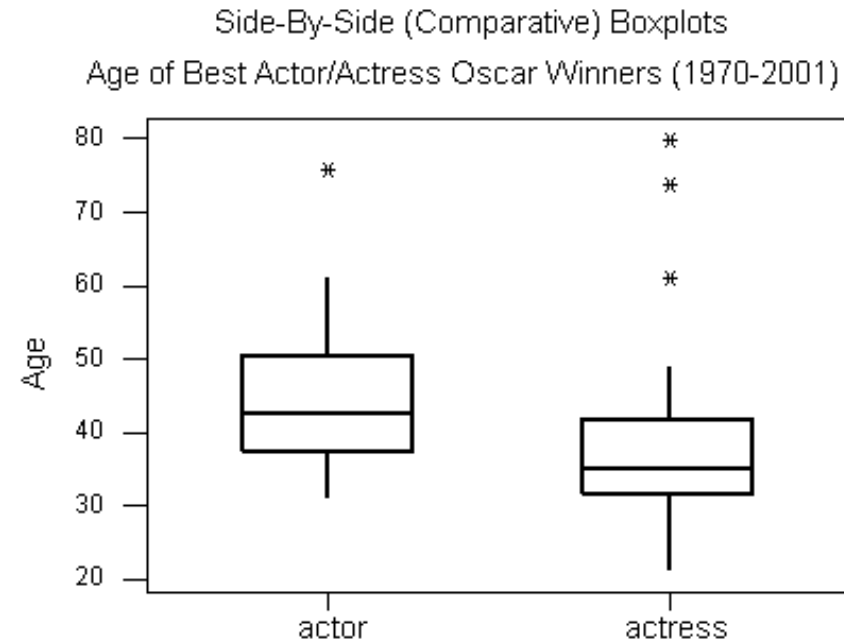
- Is the **red point** an outlier? What if we add the **blue points**?



- Red point has the **lowest z-score**.
 - In the first case it was a “**global**” outlier.
 - In this second case it’s a “**local**” outlier:
 - Within normal data range, but **far from other points**.
- It’s hard to precisely define “outliers”.
 - Can we have **outlier groups**? What about repeating patterns?

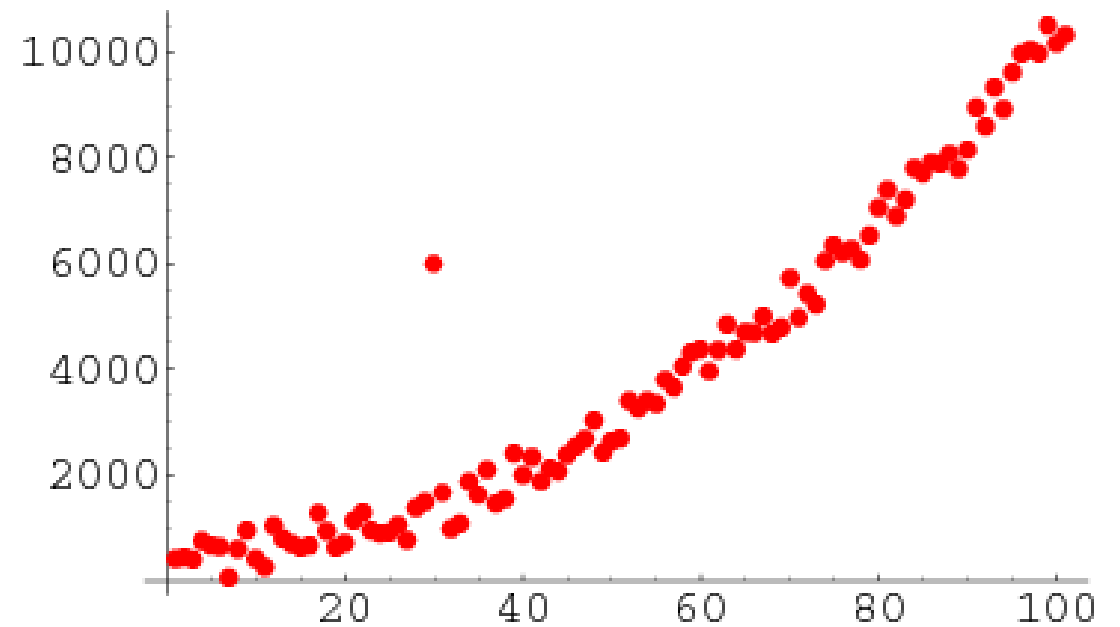
Graphical Outlier Detection

- Graphical approach to outlier detection:
 1. Look at a plot of the data.
 2. Human decides if data is an outlier.
- Examples:
 1. Box plot:
 - Visualization of quantiles/outliers.
 - Only 1 variable at a time.



Graphical Outlier Detection

- Graphical approach to outlier detection:
 1. Look at a plot of the data.
 2. Human decides if data is an outlier.
- Examples:
 1. Box plot.
 2. Scatterplot:
 - Can detect complex patterns.
 - Only 2 variables at a time.



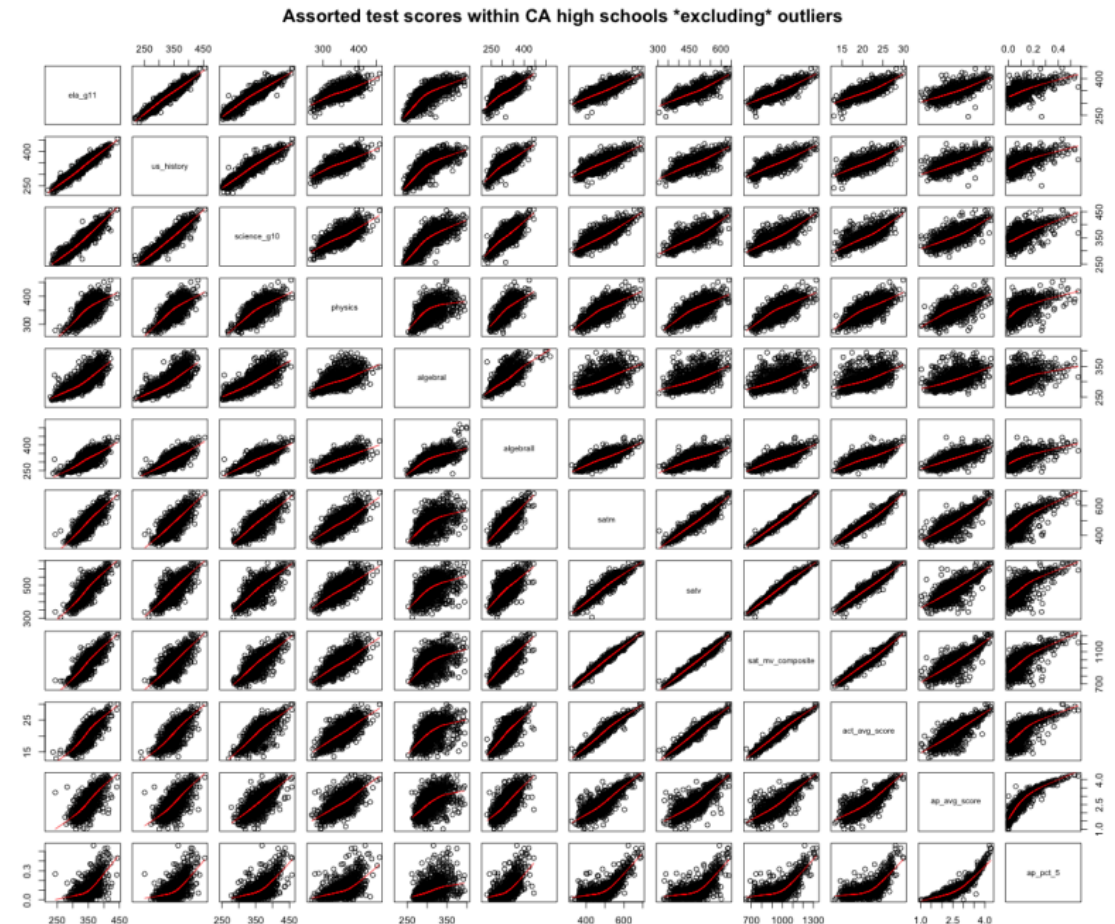
Graphical Outlier Detection

- Graphical approach to outlier detection:

1. Look at a plot of the data.
2. Human decides if data is an outlier.

- Examples:

1. Box plot.
2. Scatterplot.
3. Scatterplot array:
 - Look at all combinations of variables.
 - But laborious in high-dimensions.
 - Still only 2 variables at a time.



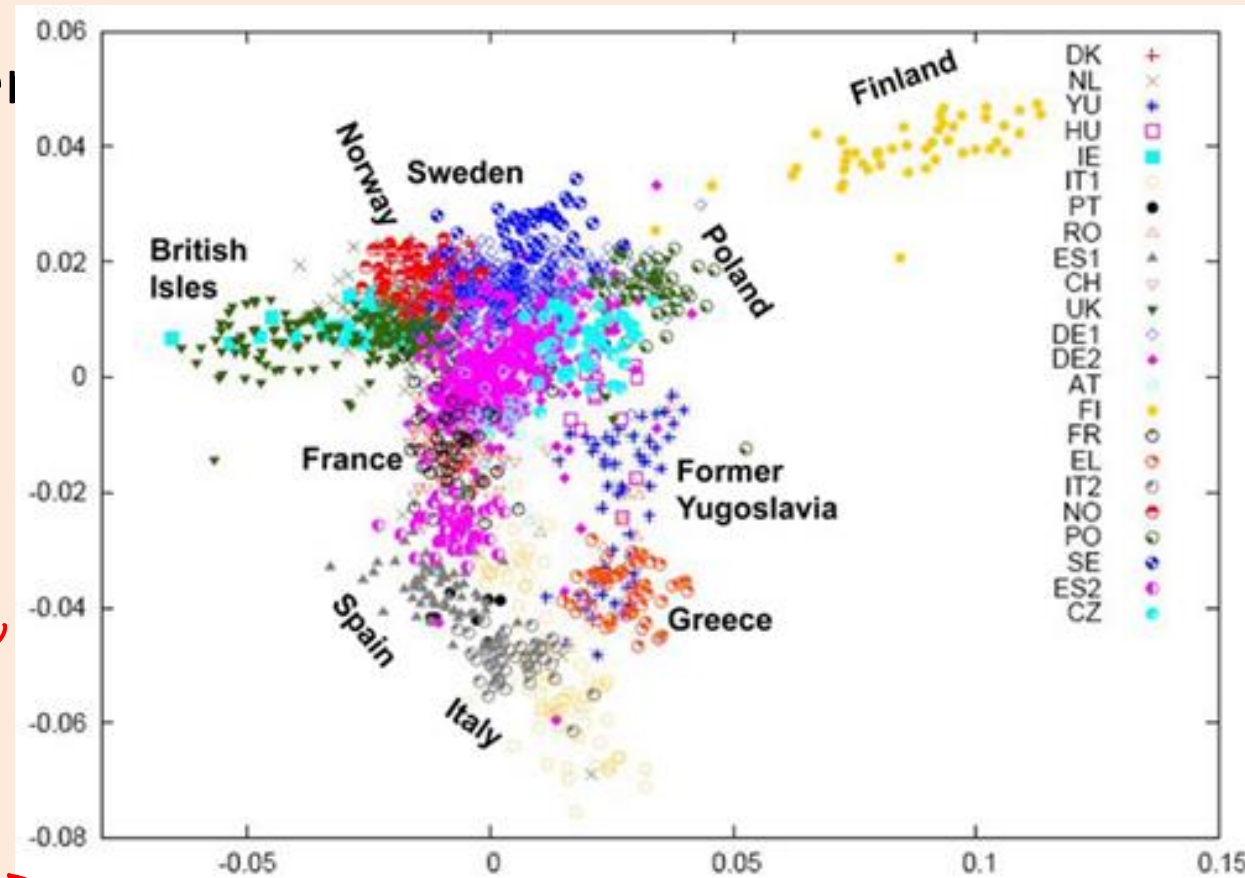
Graphical Outlier Detection

- **Graphical approach** to outlier detection:

1. Look at a plot of the data.
2. Human decides if data is an outlier

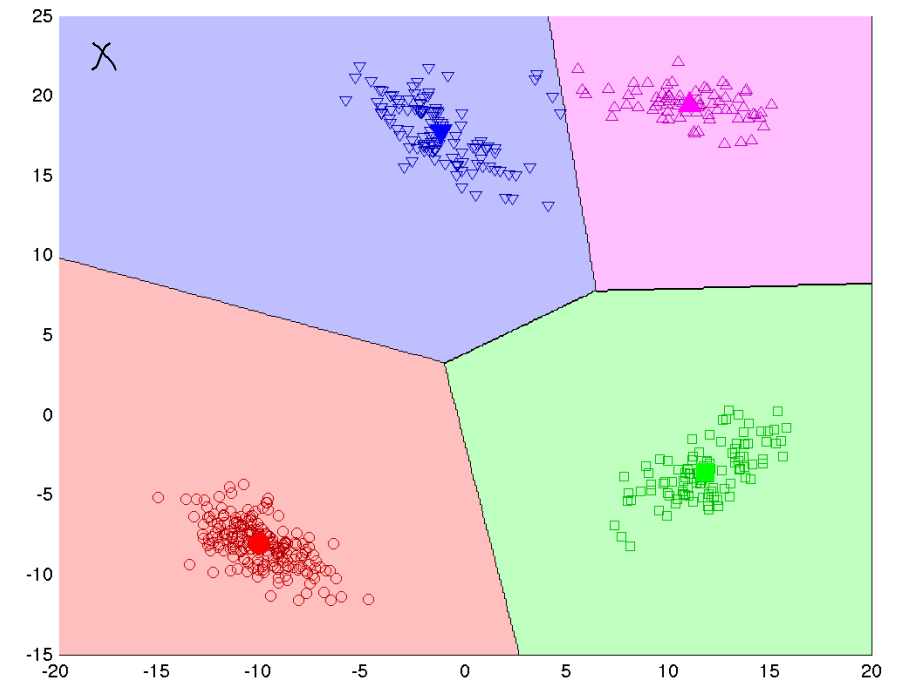
- **Examples:**

1. Box plot.
2. Scatterplot.
3. Scatterplot array.
4. **Scatterplot of 2-dimensional PCA:**
 - 'See' high-dimensional structure.
 - But **loses information** and **sensitive to outliers**.



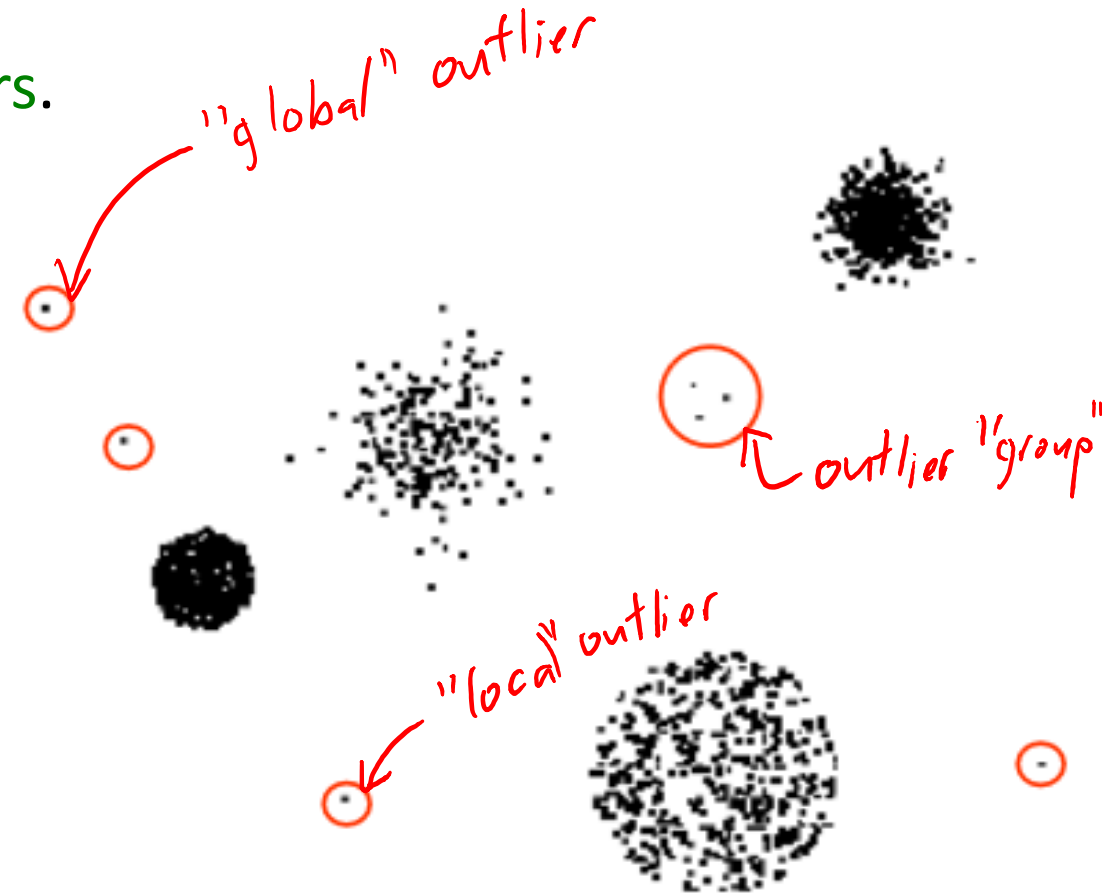
Cluster-Based Outlier Detection

- Detect outliers based on **clustering**:
 1. Cluster the data.
 2. Find **points that don't belong to clusters**.
- Examples:
 1. K-means:
 - Find points that are far away from any mean.
 - Find clusters with a small number of points.



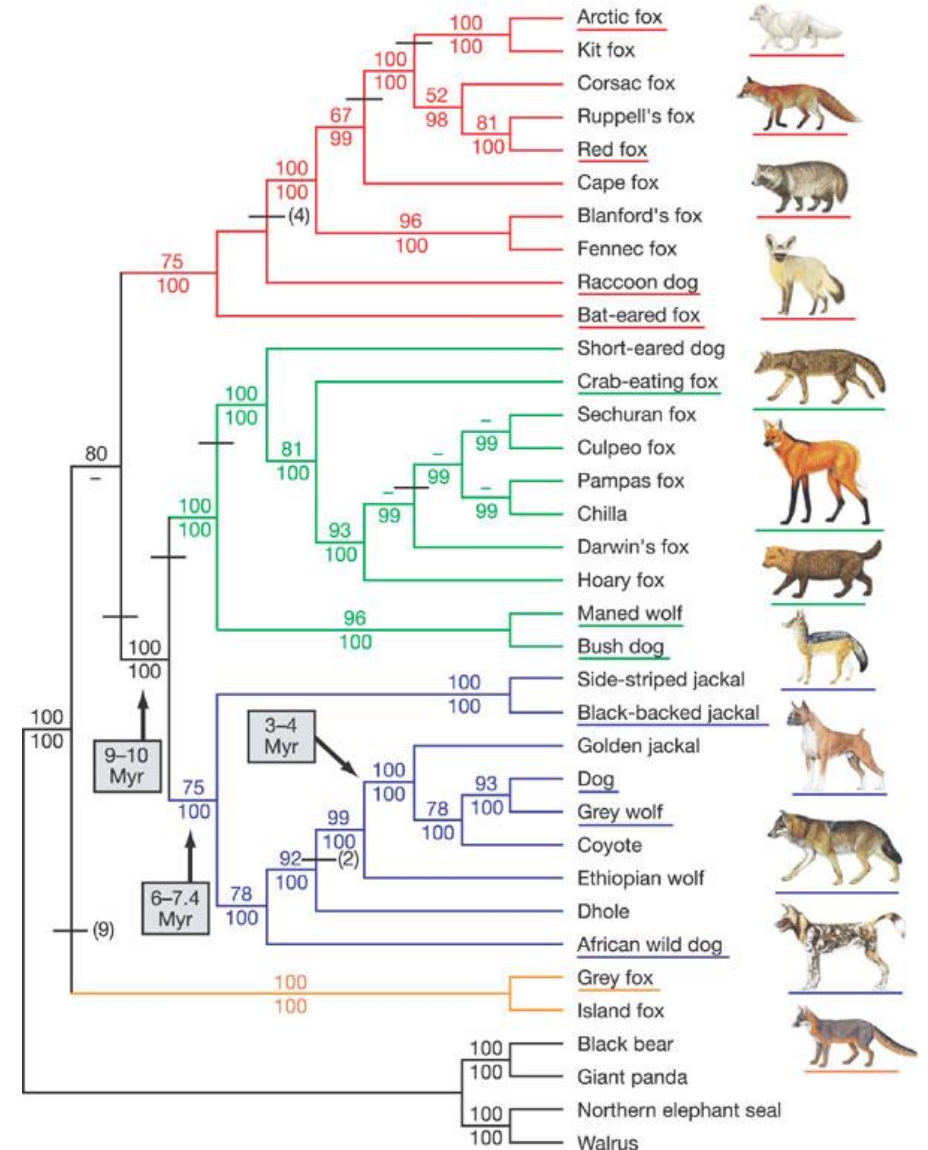
Cluster-Based Outlier Detection

- Detect outliers based on clustering:
 1. Cluster the data.
 2. Find points that don't belong to clusters.
- Examples:
 1. K-means.
 2. Density-based clustering:
 - Outliers are points not assigned to cluster.



Cluster-Based Outlier Detection

- Detect outliers based on clustering:
 1. Cluster the data.
 2. Find points that don't belong to clusters.
- Examples:
 1. K-means.
 2. Density-based clustering.
 3. Hierarchical clustering:
 - Outliers take longer to join other groups.
 - Also good for outlier groups.



Distance-Based Outlier Detection

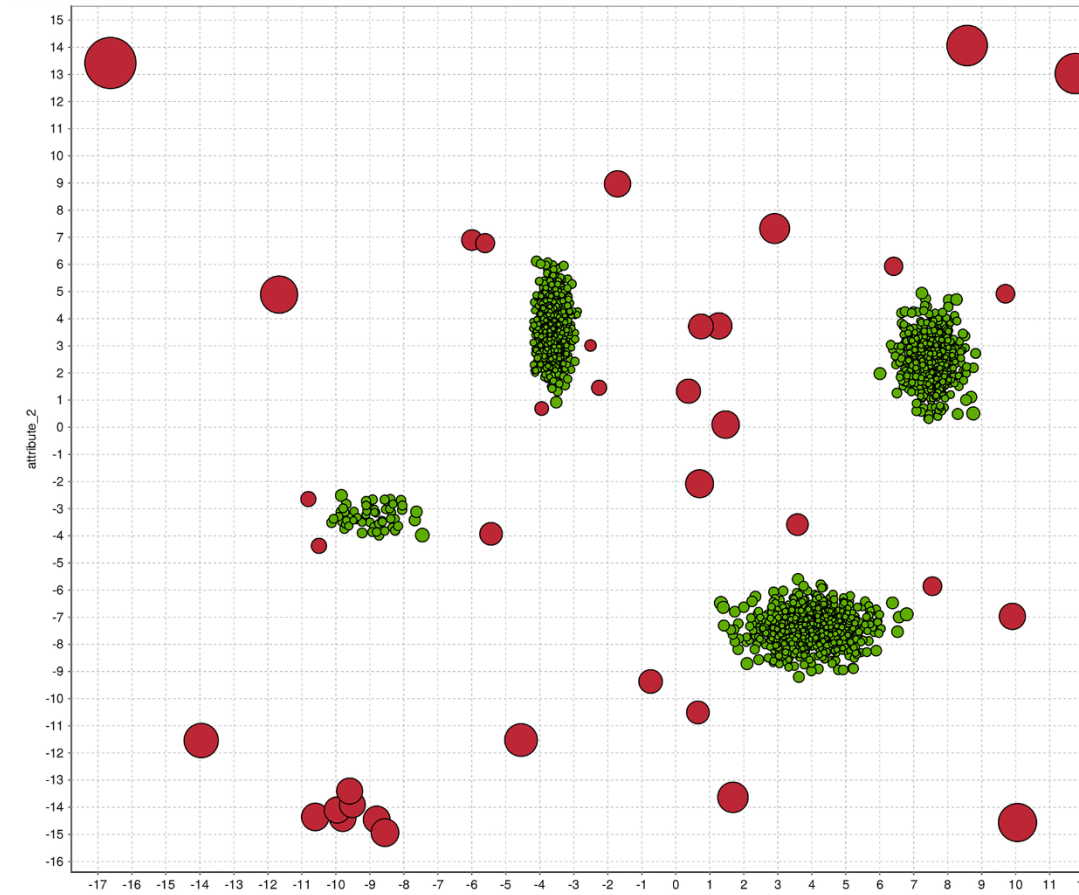
- Most outlier detection approaches are **based on distances**.
- Can we skip the model/plot/clustering and **just measure distances**?
 - How many points lie in a radius 'epsilon'?
 - What is distance to k^{th} nearest neighbour?
- UBC connection (first paper on this topic):

Algorithms for Mining Distance-Based Outliers in Large Datasets

Edwin M. Knorr and Raymond T. Ng
Department of Computer Science
University of British Columbia

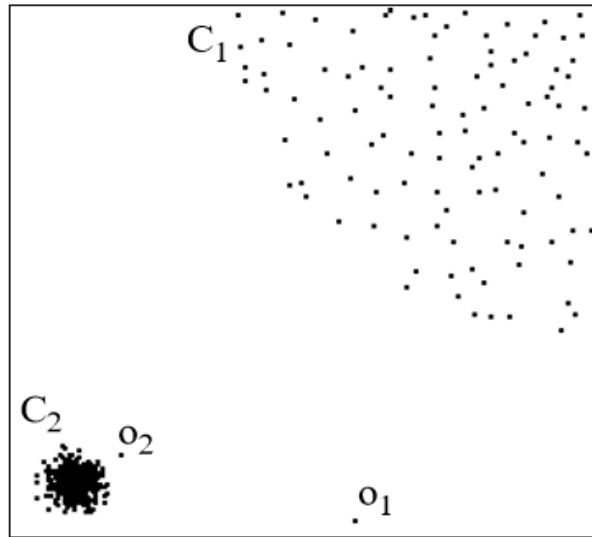
Global Distance-Based Outlier Detection: KNN

- **KNN outlier detection:**
 - For each point, compute the **average distance to its KNN**.
 - Sort the set of ‘n’ average distances.
 - Choose the biggest values as outliers.
 - **Filter out points that are far from their KNNs.**
- Goldstein and Uchida [2016]:
 - Compared 19 methods on 10 datasets.
 - **KNN best for finding “global” outliers.**
 - “Local” outliers best found with **local distance-based** methods...



Local Distance-Based Outlier Detection

- As with density-based clustering, **problem with differing densities:**



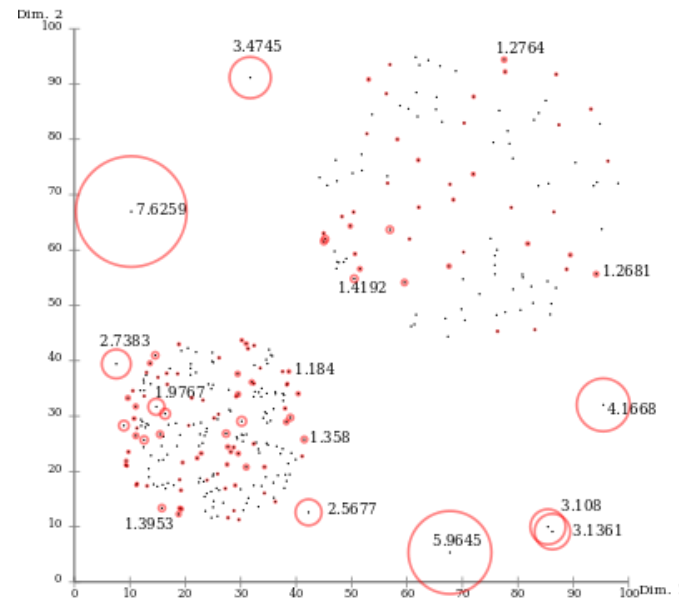
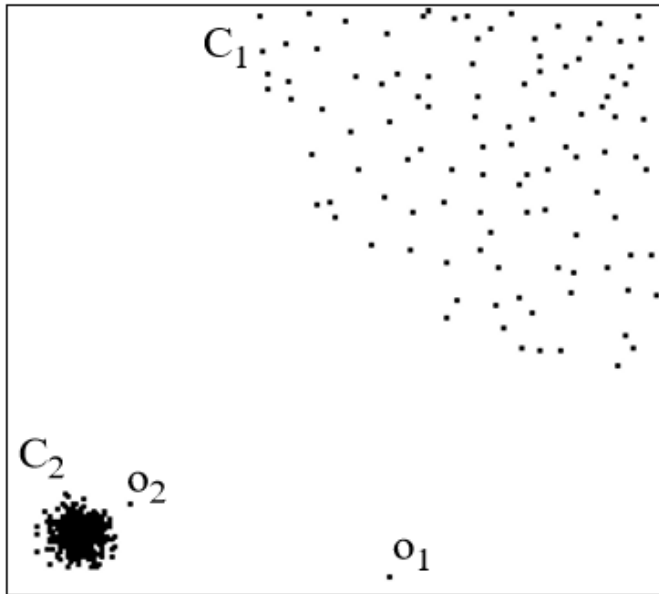
- Outlier o_2 has similar density as elements of cluster C_1 .
- Basic idea behind **local distance-based** methods:
 - Outlier o_2 is “**relatively**” **far** compared to its neighbours.

Local Distance-Based Outlier Detection

- “Outlierness” ratio of example ‘i’:

average distance of ‘i’ to its KNNs
average distance of neighbours of ‘i’ to their KNNs

- If outlierness > 1 , x_i is further away from neighbours than expected.

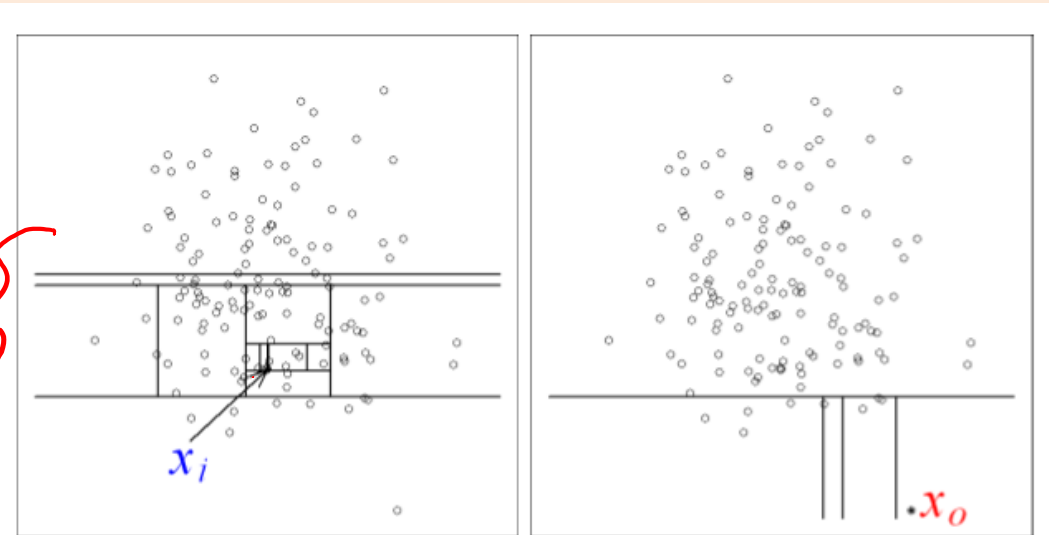


Isolation Forests

- Recent method based on random trees is **isolation forests**.
 - Grow a tree where **each stump uses a random feature and random split**.
 - Stop when each example is “isolated” (each leaf has one example).
 - The “**isolation score**” is the depth before example gets isolated.
 - Outliers should be isolated quickly, inliers should need lots of rules to isolate.

Depth 12:
– needed 12
rules to isolate
so may be inlier.

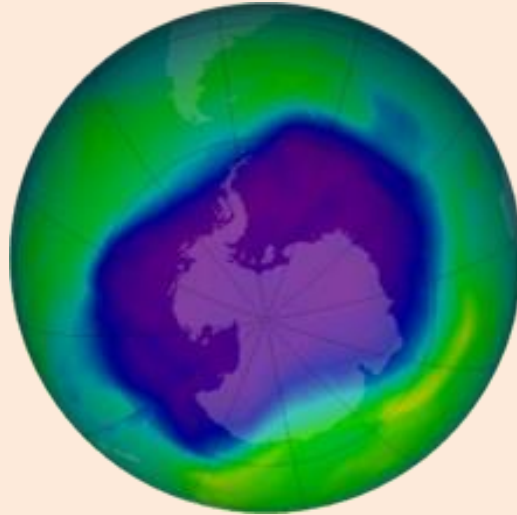
- Repeat for different random trees, take average score.



depth 4
so more likely to be outlier

Problem with Unsupervised Outlier Detection

- Why wasn't the hole in the ozone layer discovered for 9 years?



- Can be **hard to decide when to report** an outlier:
 - If **you report too many non-outliers, users will turn you off.**
 - Most antivirus programs do not use ML methods (see ["base-rate fallacy"](#))

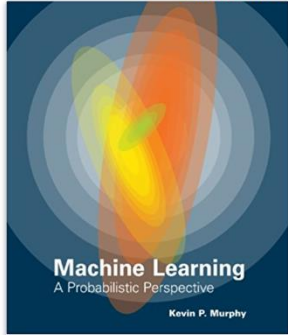
Supervised Outlier Detection

- Final approach to outlier detection is to use supervised learning:
 - $y_i = 1$ if x_i is an outlier.
 - $y_i = 0$ if x_i is a regular point.
- We can use our methods for supervised learning:
 - We can find very complicated outlier patterns.
 - Classic credit card fraud detection methods used decision trees.
- But it needs supervision:
 - We need to know what outliers look like.
 - We may not detect new “types” of outliers.

(pause)

Motivation: Product Recommendation

- A customer comes to your website looking to buy at item:



Machine Learning: A Probabilistic Perspective
Hardcover – Aug 24 2012
by Kevin P. Murphy (Author)
★★★★☆ 4 customer reviews

See all 3 formats and editions

Kindle Edition CDNS 117.34	Hardcover CDNS 123.52
-------------------------------	---------------------------------


Read with Our **Free App** 10 Used from CDNS 110.00
15 New from CDNS 99.86

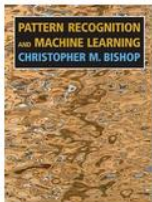

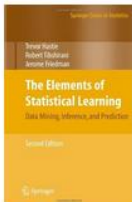
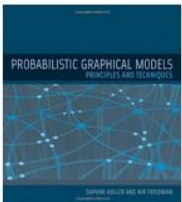
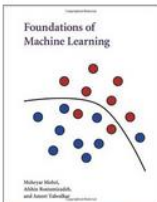
Save up to 50% on *Dummies* See more

- You want to **find similar items** that they might also buy:

Customers Who Bought This Item Also Bought

Page 1 of 20

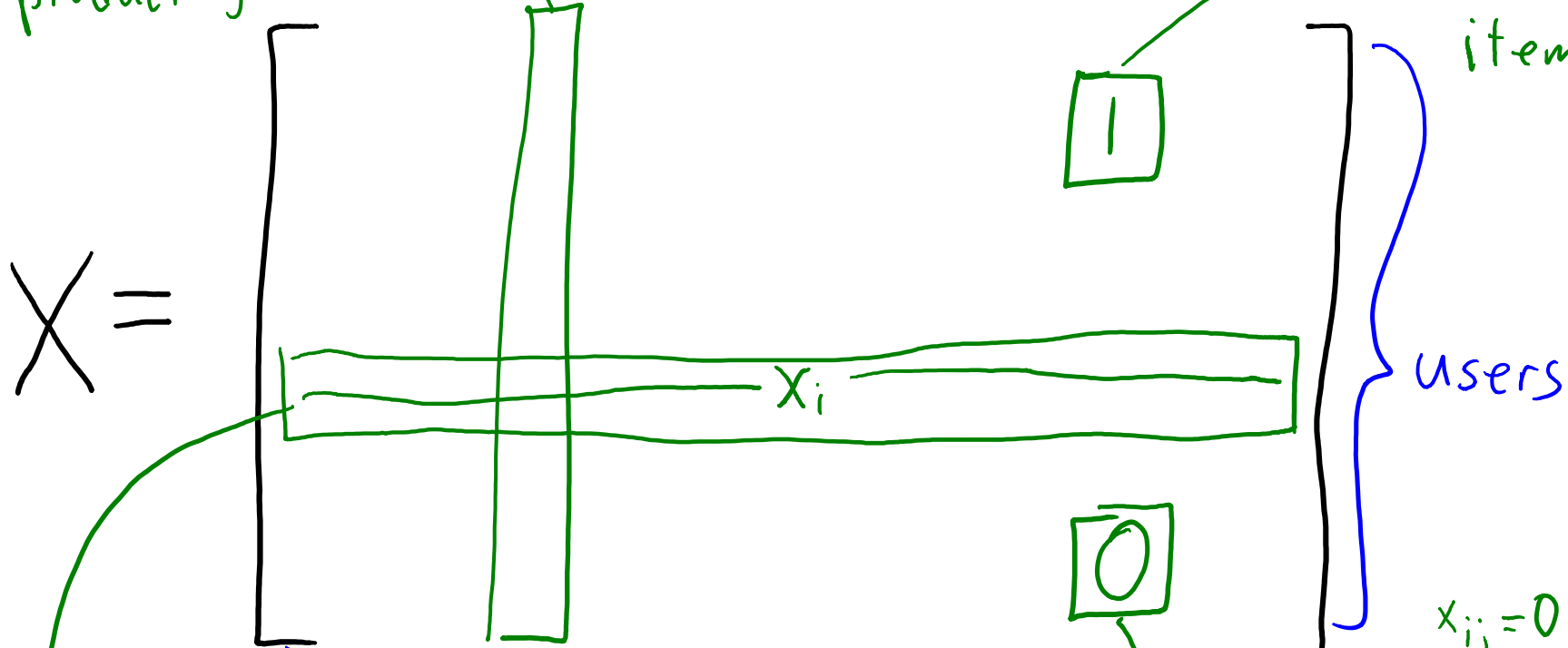


 <p>Pattern Recognition and Machine Learning (Information Science and...) Christopher Bishop ★★★★☆ 115 Hardcover \$60.76 Prime</p>	 <p>Learning From Data Yaser S. Abu-Mostafa ★★★★☆ 88 Hardcover</p>	 <p>The Elements of Statistical Learning: Data Mining, Inference, and Prediction... Trevor Hastie ★★★★☆ 50 Hardcover \$62.82 Prime</p>	 <p>Probabilistic Graphical Models: Principles and Techniques (Adaptive... Daphne Koller ★★★★☆ 28 Hardcover \$91.66 Prime</p>	 <p>Foundations of Machine Learning (Adaptive Computation and... Mehryar Mohri ★★★★☆ 8 Hardcover \$65.68 Prime</p>
--	---	--	---	--

User-Product Matrix

Column x^j gives
all users that
bought product 'j'

$x_{ij} = 1$ means
user 'i' bought
item 'j'!

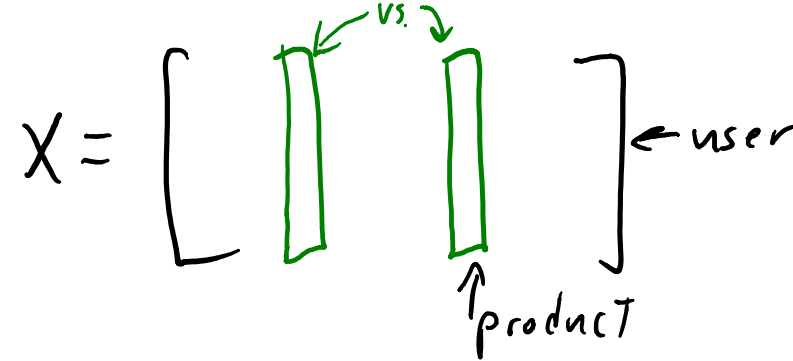


$x_{ij} = 0$ means user 'i'
has not buy item 'j'

Row x_i gives all items bought by user 'i'

Amazon Product Recommendation

- Amazon product recommendation method:



- Return the **KNNs across columns**.
 - Find 'j' values minimizing $||x^i - x^j||$.
 - **Products that were bought by similar sets of users.**
- But first **divide each column by its norm, $x^i / ||x^i||$** .
 - This is called **normalization**.
 - Reflects whether product is bought by many people or few people.

Amazon Product Recommendation

- Consider this user-item matrix:

$$X = \begin{matrix} & \text{Product 1} & \text{Product 2} & \text{Product 3} & \text{Product 4} & \text{Product 5} & \text{Product 6} \\ \text{John} & 1 & 1 & 1 & 1 & 0 & 1 \\ \text{Paul} & 1 & 0 & 1 & 0 & 1 & 0 \\ \text{George} & 1 & 0 & 1 & 0 & 1 & 1 \\ \text{Ringo} & 1 & 0 & 1 & 0 & 1 & 1 \\ \text{Yoko} & 1 & 1 & 0 & 1 & 0 & 0 \end{matrix}$$

- Product 1 is most similar to Product 3 (bought by lots of people).
- Product 2 is most similar to Product 4 (also bought by John and Yoko).
- Product 3 is **equally similar to Products 1, 5, and 6**.
 - Does not take into account that Product 1 is more popular than 5 and 6.

Amazon Product Recommendation

- Consider this user-item matrix (normalized):

$$X = \begin{matrix} & \text{Product 1} & \text{Product 2} & \text{Product 3} & \text{Product 4} & \text{Product 5} & \text{Product 6} \\ \text{John} & \frac{1}{\sqrt{5}} & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{4}} & \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{3}} \\ \text{Paul} & \frac{1}{\sqrt{5}} & 0 & \frac{1}{\sqrt{4}} & 0 & \frac{1}{\sqrt{3}} & 0 \\ \text{George} & \frac{1}{\sqrt{5}} & 0 & \frac{1}{\sqrt{4}} & 0 & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \\ \text{Ringo} & \frac{1}{\sqrt{5}} & 0 & \frac{1}{\sqrt{4}} & 0 & \frac{1}{\sqrt{3}} & \frac{1}{\sqrt{3}} \\ \text{Yoko} & \frac{1}{\sqrt{5}} & \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{2}} & 0 & 0 \end{matrix}$$

- Product 1 is most similar to Product 3 (bought by lots of people).
- Product 2 is most similar to Product 4 (also bought by John and Yoko).
- Product 3 is **most similar to Product 1**.
 - Normalization means it **prefers the popular items**.

Cost of Finding Nearest Neighbours

- With 'n' users and 'd' products, finding KNNs costs $O(nd)$.
 - Not feasible if 'n' and 'd' are in the millions.
- It's faster if the user-product matrix is sparse: $O(z)$ for z non-zeroes.
 - But 'z' is still enormous in the Amazon example.

Closest-Point Problems

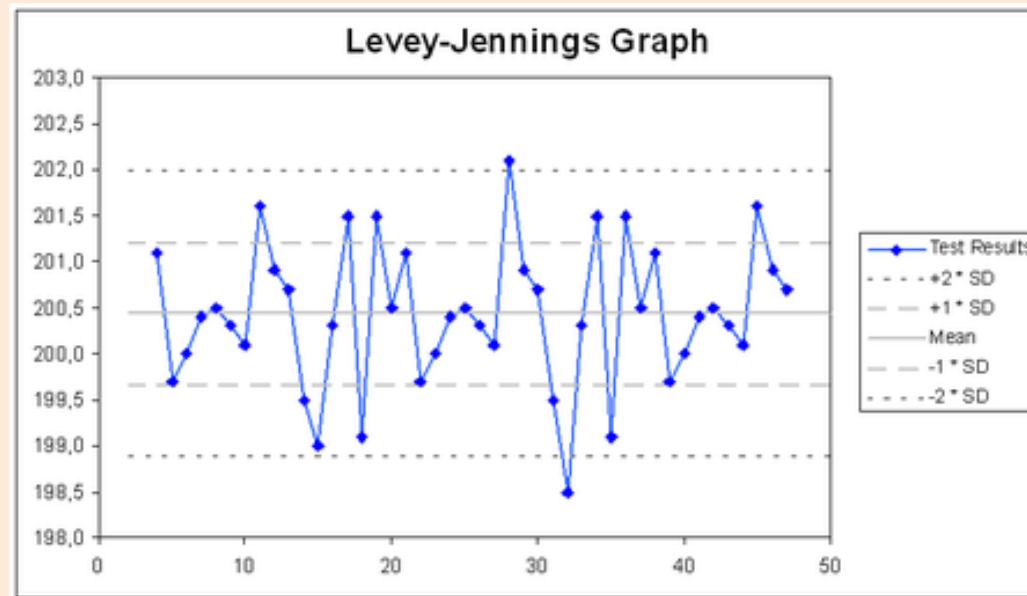
- We've seen a lot of “closest point” problems:
 - K-nearest neighbours classification.
 - K-means clustering.
 - Density-based clustering.
 - Hierarchical clustering.
 - KNN-based outlier detection.
 - Outlierness ratio.
 - Amazon product recommendation.
- How can we possibly apply these to Amazon-sized datasets?

Summary

- **Outlier detection** is task of finding unusually different example.
 - A concept that is very difficult to define.
 - **Model-based** find unlikely examples given a model of the data.
 - **Graphical** methods plot data and use human to find outliers.
 - **Cluster-based** methods check whether examples belong to clusters.
 - **Distance-based outlier detection**: measure (relative) distance to neighbours.
 - **Supervised-learning for outlier detection**: turns task into supervised learning.
- **Amazon product recommendation**:
 - Find similar items using (normalized) nearest neighbour search.
- Next time: detecting genes, viruses, plagiarism, and fingerprints.

“Quality Control”: Outlier Detection in Time-Series

- A field primarily focusing on outlier detection is **quality control**.
- One of the main tools is plotting z-score thresholds over time:



- Usually don't do tests like " $|z_i| > 3$ ", since this happens normally.
- Instead, identify problems with tests like " $|z_i| > 2$ twice in a row".

Outlierness (Symbol Definition)

- Let $N_k(x_i)$ be the **k-nearest neighbours** of x_i .
- Let $D_k(x_i)$ be the **average distance** to k-nearest neighbours:

$$D_k(x_i) = \frac{1}{k} \sum_{j \in N_k(x_i)} \|x_i - x_j\|$$

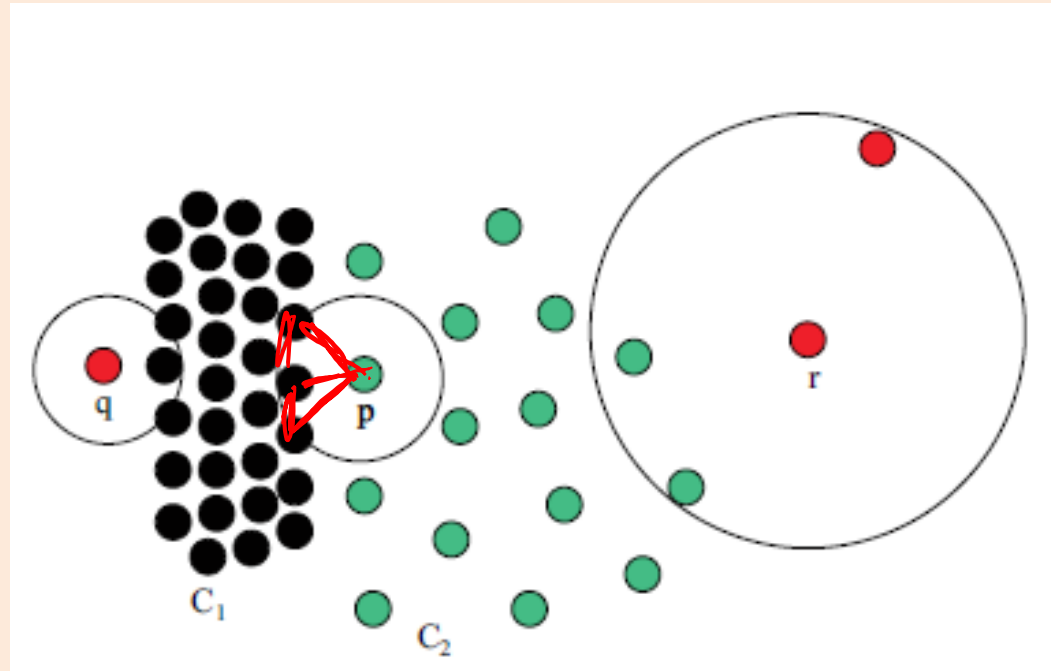
- **Outlierness** is ratio of $D_k(x_i)$ to average $D_k(x_j)$ for its neighbours 'j':

$$O_k(x_i) = \frac{D_k(x_i)}{\frac{1}{k} \sum_{j \in N_k(x_i)} D_k(x_j)}$$

- If outlierness > 1 , x_i is **further away from neighbours** than expected.

Outlierness with Close Clusters

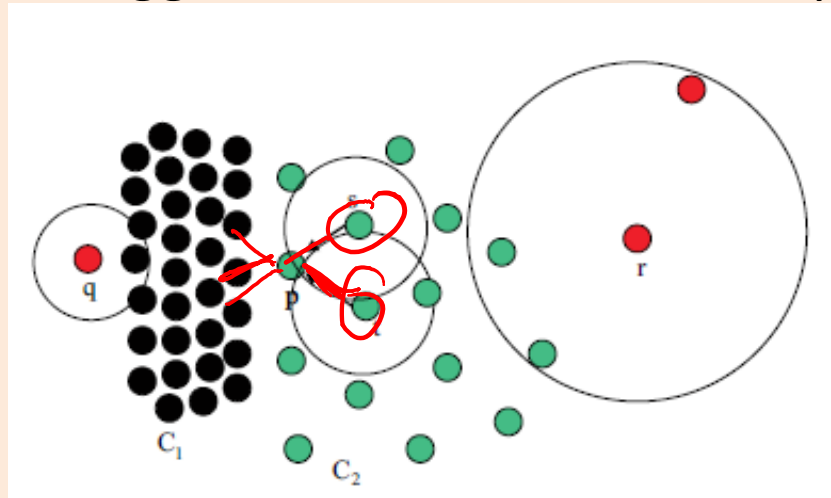
- If clusters are close, outlierness gives unintuitive results:



- In this example, 'p' has higher outlierness than 'q' and 'r':
 - The green points are not part of the KNN list of 'p' for small 'k'.

Outlierness with Close Clusters

- ‘Influenced outlierness’ (INFLO) ratio:
 - Include in denominator the ‘reverse’ k-nearest neighbours:
 - Points that have ‘p’ in KNN list.
 - Adds ‘s’ and ‘t’ from bigger cluster that includes ‘p’:



- But still has problems:
 - Dealing with hierarchical clusters.
 - Yields many false positives if you have “global” outliers.
 - Goldstein and Uchida [2016] recommend just using KNN.