Lecture 16

Outlier Analysis & Anomaly Detection

**Outlier analysis**, also known as **anomaly detection** or **anomaly analysis**, is a critical aspect of data mining and data analysis. It involves identifying data points or observations that deviate significantly from the majority of the dataset. These data points are called outliers, anomalies, or novelties, and they may indicate errors, rare events, or important insights. We touched on several aspects of identifying outliers in MTH 324 and 325. Here's an overview of outlier analysis in data mining:

*1. What Are Outliers?* *Outliers*: Outliers are data points that are significantly different from other data points in a dataset. They can be unusually high or low values, data points in different clusters, or data points with characteristics that deviate from the norm.

*2. Importance of Outlier Analysis*: Outlier analysis is important for various reasons, including: Detecting errors in data. Identifying fraud or unusual activities in finance and security. Finding unusual medical conditions or anomalies in healthcare. Discovering unexpected patterns and insights in data. Ensuring the quality and reliability of data.

*3. Techniques for Outlier Analysis*:
*Statistical Methods*: Statistical techniques like the z-score or modified z-score, Tukey's fences, and the IQR (Interquartile Range) are used to identify outliers based on statistical measures.

*Distance-Based Methods*: Distance-based methods, such as the Mahalanobis distance or Euclidean distance, can help detect outliers by measuring the distance of data points from a center point or cluster.

*Density-Based Methods*: Algorithms like DBSCAN (Density-Based Spatial Clustering of Applications with Noise) can be used to identify data points that do not fit well within dense regions.

*Clustering-Based Methods*: Outliers can be detected by analyzing the clustering results and considering data points that do not belong to any cluster as outliers.

*Machine Learning-Based Methods*: Supervised and unsupervised machine learning techniques, including Isolation Forests, One-Class SVM (Support Vector Machine), and Autoencoders, are effective for identifying outliers.

*4. Challenges in Outlier Analysis*: One of the main challenges is defining what constitutes an outlier, as it can be context-dependent. The choice of the appropriate method or technique for outlier analysis depends on the characteristics of the data and the problem domain. Handling imbalanced datasets where outliers are rare can be challenging.

*5. Visualization for Outlier Detection*: Data visualization techniques, such as scatter plots, box plots, and histograms, can help identify outliers by visually inspecting the data.

*6. Evaluation*: Evaluating the performance of an outlier detection model can be challenging, as there is often an imbalance between the number of outliers and non-outliers. Metrics like precision, recall, and F1-score can be used to assess the effectiveness of an outlier detection model.

***7. Real-World Applications***: Outlier analysis is used in various fields, including finance (fraud detection), healthcare (disease outbreak detection), manufacturing (fault detection), and more.

Outlier analysis plays a crucial role in data mining and data analysis, helping to uncover hidden insights, detect errors, and make data-driven decisions in various domains. The choice of methods and techniques depends on the specific data and problem at hand.

Resources:
1. https://statsandr.com/blog/outliers-detection-in-r/
2. https://www.digitalocean.com/community/tutorials/outlier-analysis-in-r
3. https://rpubs.com/Alema/1000582
4. https://www.geeksforgeeks.org/outlier-analysis-in-r/
5. https://www.reneshbedre.com/blog/find-outliers.html
6. https://www.r-bloggers.com/2016/12/outlier-detection-and-treatment-with-r/
7. https://universeofdatascience.com/how-to-test-for-identifying-outliers-in-r/
8. https://towardsdatascience.com/tidy-anomaly-detection-using-r-82a0c776d523
9. https://rpubs.com/michaelmallari/anomaly-detection-r
10. https://www.analyticsvidhya.com/blog/2020/12/a-case-study-to-detect-anomalies-in-time-series-using-anomalize-package-in-r/
11. https://community.sisense.com/t5/knowledge/anomaly-detection-with-sisense-using-r/ta-p/9482
12. https://www.r-bloggers.com/2018/06/anomaly-detection-in-r-2/
13. https://business-science.github.io/timetk/articles/TK08_Automatic_Anomaly_Detection.html
14. https://github.com/pridiltal/ctv-AnomalyDetection
15. https://towardsdatascience.com/tidy-anomaly-detection-using-r-82a0c776d523
16. https://rpubs.com/michaelmallari/anomaly-detection-r