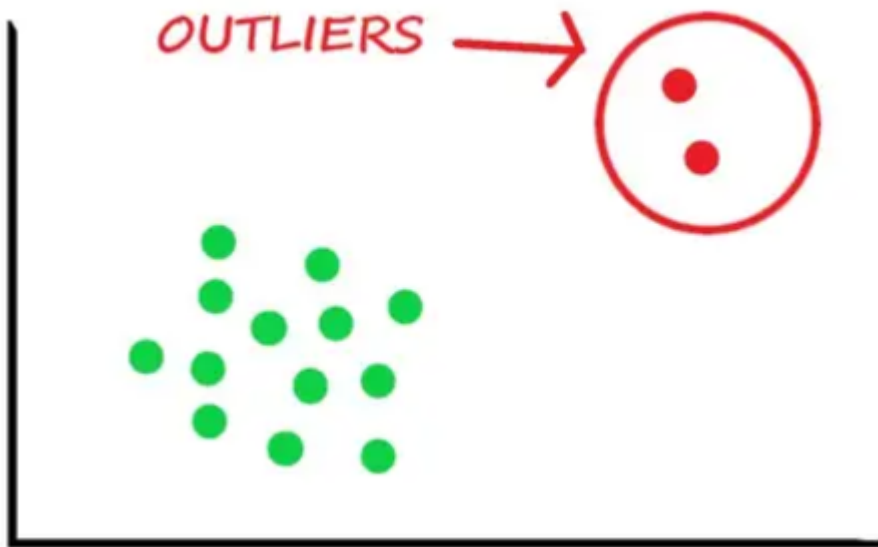


Outlier Treatment

Outliers are data points that are significantly different from the rest of the data. They are unusually high or low values that don't follow the general trend or pattern of the data set. Imagine you're looking at test scores of a class, and most of the scores are between 60 and 90, but one or two students scored 200. That 200 is an outlier because it doesn't match the rest of the scores.



Why Do Outliers Matter?

Outliers can have a big impact on statistical analysis. They can:

- **Skew averages:** For example, if you have one really high income (like \$1,000,000) in a dataset where most people earn \$40,000, the average income will be much higher than most people's actual income.
- **Affect models:** In machine learning, outliers can lead to incorrect predictions because the model may try to "fit" the data too closely, including the outlier, and end up being less accurate.

How to Identify Outliers?

There are several ways to spot outliers, and one of the most common methods is through **visualization** or **statistical rules**. Let's break these down:

1. Visual Methods:

- **Box Plot:** This is a simple chart that shows the distribution of data. Outliers are typically marked as dots outside the "whiskers" (the lines of the box plot). If the data points fall far outside the normal range (1.5 times the interquartile range), they are considered outliers.
- **Scatter Plot:** A scatter plot shows individual data points. If one point is far away from the rest, it could be an outlier.

For example

In a **scatter plot** of a city's housing prices, most of the data might be clustered around \$200,000 to \$500,000, but a few points might show properties worth \$5 million. These high-priced homes are the outliers.

2. Statistical Methods:

- **Z-scores:** This is a common statistical measure that tells us how many standard deviations a data point is from the mean. If a Z-score is greater than 3 or less than -3, it's considered an outlier.
- **IQR (Interquartile Range):** The IQR is the range between the first and third quartile of the data. Outliers are typically any data points that fall outside of 1.5 times the IQR.

How to deal with Outliers?

- **When to use it:** If the outlier is a mistake (like a typo in data entry), or if it's so extreme that it's not relevant to your analysis, removing it can make your data set more accurate.

Example

- If you're analyzing the weight of people in a fitness study, and one person's weight is recorded as 500 kg when everyone else is in the range of 50-100 kg, it's likely a data entry error. You would remove this outlier.

2. Transform the Data:

- **When to use it:** If the outliers are valid but distort the analysis, you can apply transformations like taking the **log** of the data to reduce the effect of outliers.

Example

- In a dataset of house prices, if most houses are worth between \$100,000 and \$300,000 but a few are worth over \$10 million, applying a log transformation can compress the scale and make the data more manageable.

3. Cap or Floor the Outliers (Winsorizing):

- **When to use it:** This is when you replace the extreme outliers with a more reasonable value, such as the maximum or minimum value within an acceptable range.

Example

- If a salary dataset has a salary of \$1 million but most salaries are between \$30,000 and \$100,000, you might decide to cap the maximum salary to \$200,000. This way, you limit the effect of extreme values without completely removing them.

4. Use Robust Methods:

- **When to use it:** Sometimes, you don't want to remove outliers, especially if they contain valuable information (like rare events). In this case, you can use models that are **robust** to outliers.

Example

- In regression analysis, you might use **robust regression** techniques that don't let outliers heavily influence the results.

Few Real-life Examples of Outliers:

1. **Real Estate Prices:** Imagine you're analyzing the prices of homes in a city. Most homes are priced between \$100,000 and \$500,000, but there are a few that are priced at \$10 million. These high-priced homes are outliers. They could either be ultra-luxury properties or errors in data entry. How you treat these outliers depends on the context.
2. **Student Exam Scores:** If most students in a class score between 60 and 90 on a test, but one student scores 200, this score is an outlier. You might need to check if it was an error, or if the student had special circumstances. If it's a real score, you might decide to keep it, but note that it could influence class averages.
3. **Customer Transactions:** If you're looking at the amount spent by customers in an online store, most transactions could be in the range of \$20-\$200, but one customer

might have spent \$10,000 on a luxury product. This is an outlier. Depending on the business, this could be a valid transaction or something to investigate further.

Conclusion:

Outliers are data points that stand out from the rest because they are much higher or lower than the others in a data set. They can be caused by errors, rare events, or legitimate variation. Dealing with outliers depends on their cause and the context of the data. You can remove, transform, or leave them, or use robust methods that minimize their influence. Identifying and understanding outliers is crucial for making better decisions based on data!