

What is k-Means Clustering?

k-means is a centroid-based partitioning method used to divide a dataset into k clusters. The idea is to group similar data points together and find underlying patterns or structure in the data. Each cluster is represented by a centroid (a central point), which is the mean of all data points in that cluster.

12 Inputs

- D: A dataset containing n objects in Euclidean space
- k: The number of clusters to form

▲ Goal of k-Means

To divide p into k clusters c1, c2, ..., ck such that:

- Each object is assigned to exactly one cluster
- The objects in the same cluster are similar
- Objects in different clusters are dissimilar

This is achieved by minimizing the sum of squared errors (SSE) within clusters (also called withincluster variation).

✓ Objective Function: Sum of Squared Errors (SSE)

$$E = \sum_{i=1}^{k} \sum_{p \in C_i} \operatorname{dist}(p, c_i)^2$$

Where:

- p is a point in cluster Ci
- ci is the centroid (mean) of cluster ci
- dist(p, ci) is the Euclidean distance between p and ci

Goal: Make each cluster compact (low internal variation) and distinct (low overlap).



How the k-Means Algorithm Works

Step-by-step Algorithm:

- 1. Initialize: Randomly choose k data points as the initial centroids.
- 2. **Assign**: For each point, assign it to the cluster with the **nearest centroid** (using Euclidean distance).
- 3. **Update**: Recalculate the centroid of each cluster as the **mean** of all data points in that cluster.
- 4. Repeat Steps 2–3 until:
 - No changes in cluster assignments OR
 - A maximum number of iterations is reached

This process is called **iterative relocation**.

II Example of k-Means Clustering

Let's say we have 2D points (as shown in a figure), and we want to cluster them into k = 3 groups:

- 1. Start with 3 random points as centroids (marked with +)
- 2. Assign nearby points to the nearest centroid
- 3. Update the centroids to the mean of the assigned points
- 4. Repeat until the centroids no longer move significantly

Important Observations

- **Not guaranteed** to find the global optimal solution. May get stuck in a **local optimum**, depending on the initial centroids.
- To improve results:
 - Run the algorithm multiple times with different random seeds
 - Use smarter initialization methods (e.g., k-means++)

Time Complexity

 $O(n \cdot k \cdot t)$

Where:

- n : number of objects
- k : number of clusters
- t : number of iterations (usually small, <100)

Efficient for large datasets but still sensitive to dimensionality and initial values.

Variants & Extensions

1. k-Modes

- Used for nominal (categorical) data
- Replaces the mean with the **mode** (most frequent value)
- Uses dissimilarity measures appropriate for categorical data

2. k-Prototypes

- Handles mixed data types (numeric + nominal)
- Combines both mean and mode techniques

▲ Limitations of k-Means

- 1. User must choose k in advance
 - Hard to know the correct value of k
 - Solution: Try a range of k values and choose based on evaluation criteria (like Elbow method or Silhouette score)
- 2. Sensitive to outliers and noise
 - One outlier can shift the mean a lot
- 3. Not good for clusters with:
 - Non-convex shapes
 - Different densities or sizes
- 4. Mean must be defined
 - Not suitable for datasets where mean is not a valid central measure (e.g., strings or purely categorical data)



Making k-Means More Scalable

Several techniques help improve performance on large datasets:

1. Sampling

Run k-means on a representative subset

2. Filtering

Use spatial indexes (like KD-Trees) to reduce computation

3. Microclustering

- Pre-cluster objects into microclusters (tiny groups)
- Then apply k-means to these microclusters instead of the original data
- Used in **BIRCH** algorithm (discussed later in your book)

Summary: Advantages & Disadvantages

~	Advantages	>	K	Disadvantages

Simple and easy to implement Need to choose k manually

Efficient and fast Sensitive to initial seeds

Works well on spherical clusters Poor with non-convex shapes

Sensitive to noise and outliers Scalable to large datasets

Can be adapted for mixed data May converge to a local minimum

★ 1. What Are Outliers?

An outlier is a data object that is significantly different from other data points in a dataset.

It's as if the outlier was generated by a different process or mechanism than the rest of the data.

Main Idea

Outliers do not follow the same distribution or pattern as the majority of data. Therefore, they are considered abnormal, while the rest of the data is called normal or expected.

Example 12.1: Gaussian Distribution with Outliers

Imagine a set of data points that follow a Gaussian (normal) distribution — forming a bell-shaped curve. Now suppose there are a few data points far away from this main group, in a region marked R (as shown in a figure in the book).

These points in Region R:

- Do **not fit** the main pattern
- Are statistically distant
- Are likely to be **outliers**

Outliers vs. Noise

Outliers

- Systematic deviation from expected behavior
- May be generated by different processes
- Interesting and valuable to detect
- Example: Credit card fraud, unusual network activity, rare diseases

Noise

- Random errors or fluctuations in data
- Usually caused by sensor error, recording mistake, etc.

- Not meaningful, should be removed
- Example: Typo in a dataset, faulty sensor reading

Important Distinction:

Noise is random and usually not actionable.

Outliers are non-random and may indicate important patterns or problems.

Credit Card Fraud Example

A customer usually buys small meals and coffee. One day, they:

- Buy a large lunch (normal variation)
- Have two coffees (slightly odd, but still within normal limits)

These are **not outliers**, just **noise** or **natural variability**.

But if that same customer:

Buys five expensive electronics at midnight from a foreign country

That may be an outlier, suggesting fraudulent activity.

- ▲ Flagging normal variation as outliers leads to:
- High false positives
- Customer dissatisfaction
- Increased verification costs

So, it's crucial to **filter out noise** before detecting outliers.

Why Are Outliers Interesting?

- Outliers are suspicious or special
- They may indicate:
 - Fraud
 - Faulty equipment
 - Malicious activity
 - Scientific discoveries
- Therefore, detecting them can be valuable in real-world applications

*** How Are Outliers Detected?**

By making assumptions about the "normal" data, such as:

- It follows a specific distribution (e.g., normal distribution)
- It has a particular pattern or density

We then check which data points violate these assumptions significantly.

If a data point breaks the expected pattern, it's flagged as a potential outlier.

Q Outlier Detection vs. Novelty Detection

Feature	Outlier Detection	Novelty Detection
Goal	Find rare, abnormal events	Detect new, previously unseen patterns
Context	Often assumes data is static	Works in evolving or streaming data
Behavior	Once outlier, always outlier	New patterns become normal over time
Example	Fraudulent transaction	Emerging topics in social media

Example:

- A new social media topic first appears and looks like an outlier.
- Over time, as more people discuss it, it becomes part of the **new normal**.
- Novelty detection identifies and tracks this change.

⊗ Summary: Key Takeaways

- Outliers are data points that are significantly different from the majority.
- They are **not the same** as **noise**, which is just **random error**.
- Outliers are useful and can indicate important events or patterns.
- Detecting them involves understanding the structure of the normal data.
- Outlier detection and novelty detection are related, but **not the same**.
- Noise should be removed first to prevent misidentification.

Would you like to continue with:

- Types of outliers (global, contextual, collective)?
- Outlier detection techniques (statistical, clustering-based, etc.)?
- Or an example of how to detect outliers in a real dataset using Python or a tool like Excel?

Let me know how you'd like to proceed!



12.1.2 Types of Outliers – Full Explanation

Outliers come in three main types:

- 1. Global Outliers
- 2. Contextual (or Conditional) Outliers
- 3. Collective Outliers

Global Outliers (aka Point Anomalies)

Definition:

A global outlier is a data object that deviates significantly from the rest of the dataset — regardless of context.

These are the simplest and most commonly studied type of outliers.

Example 12.2:

Recall **Figure 12.1** (from the previous section). Most points follow a normal (Gaussian) distribution, but the points in **Region R** are far from the rest.

These are **global outliers** because they are **isolated** and **very different** from the general data.

Key Point:

Detecting global outliers requires a good measure of deviation. This could be:

- Distance from a central value (like the mean or median)
- Statistical probability (e.g., if a value is beyond 3 standard deviations)
- Clustering density, etc.

These measurements depend on the application domain.

© Use Cases:

- Intrusion detection: A computer suddenly sending thousands of data packets might be a global outlier.
- Financial audits: A transaction that doesn't follow regulations is flagged as a global outlier.

Contextual Outliers (aka Conditional Outliers)

✓ Definition:

A contextual outlier is a data object that is only considered an outlier within a specific context (like time, location, or other conditions).

Same value might be normal in one context, and abnormal in another.

Service Example:

"The temperature today is 28°C. Is it an outlier?"

- In Toronto during winter: YES, it's an outlier.
- In Toronto during summer: NO, it's normal.
- *context matters. So, context matters.*

How it works:

We divide the attributes into two types:

Туре	Description	Example
Contextual Attributes	Define the situation or setting	Date, Location
Behavioral Attributes	Define the values being analyzed	Temperature, Pressure, Humidity

Rule:

A value is an outlier **only within the context** of its contextual attributes.

For example, a customer spending 90% of their credit limit might be fine for low-income customers, but **abnormal** for high-income customers.

Local vs. Contextual:

- A **local outlier** (from density-based methods) is a **simplified version** of a contextual outlier detected based on local density, not domain-defined context.
- Global outlier detection is just a special case of contextual detection where no context is considered (i.e., the context is the whole dataset).

© Use Cases:

- Credit card fraud detection (Example 12.3):
 - High balance may not be an outlier for one customer group, but may be for another.
 - Detecting these can lead to business opportunities like increasing credit limits.

▲ Challenges:

- Requires domain knowledge to define meaningful contexts.
- Not always easy to get high-quality contextual data.
- Simple groupings (e.g., GROUP BY in SQL) might lead to insufficient data or noise.

✓ Better Method:

Use proximity in contextual space to form context dynamically — discussed later in Section 12.4.

Collective Outliers

Definition:

A collective outlier is a group (subset) of data objects that are normal individually, but abnormal as a group.

Group behavior is what makes them suspicious — not the individual points.

Example (Supply Chain):

- One delayed order: normal
- 100 delayed orders on the same day: collective outlier

Even if each order isn't abnormal alone, the group behavior signals a problem.

Example 12.4:

In Figure 12.2, black dots are densely clustered together, forming a collective outlier. Each black point by itself looks normal, but together they form an unusual pattern.

© Use Cases:

- **Cybersecurity**: A single computer sending denial-of-service (DoS) packets isn't an issue. But **multiple computers** doing this together signals a **coordinated attack**.
- Stock market fraud: A few trades are normal. But many trades of the same stock between the same parties in a short time? Possibly market manipulation.

Key Requirements:

- Need to understand relationships between data objects:
 - Distance
 - Similarity
 - Sequence
- This makes collective outlier detection more **complex** than global or contextual detection.

Summary Table

Туре	Outlier Based On	Example	Notes
Global Outlier	Deviation from entire dataset	A huge money transfer vs. small regular ones	Simplest type
Contextual Outlier	Deviation within a context	28°C in Toronto winter	Needs context attributes
Collective Outlier	Abnormal group behavior	100 shipments delayed in one day	Group-based, not individual

Final Thoughts

- A data set can contain multiple types of outliers at once.
- A data object can be both a **global** and a **contextual outlier**, for example.
- Outlier types should be chosen based on the application:
 - Fraud detection
 - Network monitoring
 - Marketing insights
- Detection techniques differ and rely on:
 - Statistical assumptions
 - Domain knowledge
 - Data structure (individual vs. group)

Would you like me to move on to the next section (e.g., 12.1.3 or 12.2), or would you like a diagram/summary table of all 3 types for revision?

12.1.3 Challenges of Outlier Detection – Full Explanation

Outlier detection is a **powerful tool** used in many domains, such as fraud detection, cybersecurity, finance, healthcare, and more.

However, **detecting outliers is not easy**. There are several significant challenges that make this task complex. Let's go through each of them one by one.

1 Modeling Normal Objects and Outliers Effectively

Q Challenge:

Outlier detection depends heavily on the ability to model what is "normal" — and by contrast, what is "abnormal."

But what exactly is *normal behavior*? In many real-world applications, it's hard to define.

⚠ Why it's difficult:

- Normal behavior can vary widely across different scenarios.
- It's almost impossible to list all possible "normal" cases.
- Sometimes, behavior that appears unusual may still be valid under certain rare conditions.

As a result, creating a perfect model of normality is extremely hard.

Gray areas:

The boundary between **normal** and **outlier** is not always clear-cut — there's often a **gray zone** where it's **subjective** whether a data point is truly an outlier.

✓ Solutions:

There are two types of approaches:

- Binary classification: Each data point is marked as either "normal" or "outlier."
- Score-based methods: Each data point gets an "outlier score" indicating how strongly it deviates from normal behavior.

Score-based methods are often more flexible because they handle the gray areas better.

2 Application-Specific Outlier Detection

Q Challenge:

There is **no one-size-fits-all** solution for outlier detection.

Different applications:

- Have different types of data
- Need different distance/similarity measures
- Interpret deviations differently

Example 1: Clinical Data

In healthcare, even small deviations in measurements (like blood pressure or heart rate) can be critical and may indicate a serious condition. So:

A small deviation = Outlier

II Example 2: Marketing Data

In marketing, user behavior fluctuates a lot — so:

A larger deviation may be needed to classify something as an outlier

★ Key Point:

- The choice of how you **model the data** and what counts as an outlier depends entirely on the **domain**.
- Therefore, customized methods are needed for each application.

It's impractical to create a universal outlier detection algorithm that works well everywhere.

3 Handling Noise in Outlier Detection

Q Challenge:

Noise is a natural part of most real-world datasets — and it can confuse outlier detection methods.

What is noise?

- Random errors or irrelevant data
- May appear as slightly wrong values, or
- Missing data fields

Why it's a problem:

- Noise may look like outliers (but it's not meaningful)
- True outliers may be hidden among noisy data
- Some outliers might be disguised as noise and not detected

Result:

- False positives: Noise is misclassified as an outlier
- False negatives: True outliers are missed because they are buried in noisy data

Solution:

- Improve data preprocessing to clean noise and impute missing values
- Use robust outlier detection algorithms that can handle noise

Understandability (Interpretability)

Q Challenge:

In some situations, it's not enough to say, "This is an outlier." We also want to **understand why** it's an outlier.

Why it matters:

- In business, healthcare, and legal applications, decisions must be justified.
- Users want to know:
 - What made this data point an outlier?
 - Which attributes were responsible?

Example:

Statistical methods can provide reasoning:

"This data point has a **low probability** of being generated by the same statistical process that produced the rest of the data — so it's an outlier."

The lower the probability, the stronger the case for being an outlier.

Summary of Key Requirements for Interpretability:

- Clear explanation of why a point is flagged
- Breakdown of contributing factors (attributes)
- Quantitative scores or probabilities

Summary Table of Challenges

Challenge Explanation

Modeling Normal & Outliers Hard to define "normal" behavior; gray zones make binary labeling tricky

Challenge Explanation

Application-Specific Detection Different domains need different methods, metrics, and

thresholds

Handling Noise Noise can mimic or hide outliers, lowering detection

accuracy

Understandability Users need to know why something is an outlier;

justification is crucial

Final Thoughts

- Outlier detection is essential, but it is not trivial.
- Each of the above challenges must be addressed to build **reliable**, **robust**, and **interpretable** outlier detection systems.
- Future sections in this chapter (like 12.2: Outlier Detection Methods) will explore how different methods try to overcome these challenges using:
 - Statistical approaches
 - Distance-based techniques
 - Clustering, density, classification, and more

Would you like me to explain Section 12.2: Outlier Detection Methods next? Or provide diagrams or real-world case studies to go with this summary?

12.2 Outlier Detection Methods – Full Explanation

Outlier detection methods vary widely, and different approaches are used depending on:

- 1. Whether labeled data is available.
- 2. What assumptions are made about how normal and outlier data behave.

To organize the field, we classify methods in two orthogonal (independent) ways:

1. Based on Label Availability

- Supervised methods: Require labeled data (both normal and outliers).
- Semi-supervised methods: Use a small set of labeled data (usually only normals or only outliers).
- Unsupervised methods: No labeled data required.

◆ 2. Based on Assumptions About Normality

Some methods assume that:

- Normal data is clustered or patterned.
- Outliers are **few** and behave **differently** than the rest.

Let's now explore each method in detail.

12.2.1 Supervised, Semi-Supervised, and Unsupervised Methods

Supervised Methods

Definition: Supervised outlier detection requires **expert-labeled data**, where we know which instances are normal and which are outliers.

The task becomes a classification problem:

- Use the labeled sample to train a model.
- Predict new data as either normal or outlier.

How it works:

- Train a classifier using labeled data (e.g., decision trees, SVM, neural networks).
- Apply the classifier to predict if a new object is an outlier.

★ Variants:

- Sometimes only normal data is labeled: Build a model of normal behavior → Anything that doesn't fit = outlier.
- Other times **outliers** are labeled: Build a model of outliers → Everything else is considered normal.

Challenges:

- 1. Class imbalance:
 - Outliers are rare compared to normal data.
 - Classifiers may learn to **ignore outliers** because they're so few.
 - Solution: Use techniques like **oversampling** outliers (duplicate them in training).
- 2. Lack of outlier representation:
 - Even with labels, rare outliers may not be well represented in the training data.
 - Poor generalization to new types of outliers.
 - Solution: Some methods generate synthetic outliers to simulate diverse behavior.
- 3. Focus on recall (sensitivity):
 - It's more important to catch all outliers than to perfectly label normals.
 - In fraud detection, missing one fraud (false negative) is more dangerous than incorrectly flagging a few normals (false positives).

✓ Summary:

Strength: High accuracy when good labeled data is available.

• Weakness: Hard to obtain enough labeled outliers; models can be biased by imbalance.

Real-life examples:

- Spam detection (labeled spam vs. legitimate emails).
- Credit card fraud detection (known fraudulent transactions).

Unsupervised Methods

Definition: These methods are used when **no labels** are available. You don't know which data points are normal or outliers.

Key assumption:

Normal data follows a pattern (e.g., clustered), while outliers:

- Are few.
- Do not fit any pattern.
- Are far from other data points in the feature space.

How it works:

- Use clustering, density, distance, or statistical distribution to identify anomalies.
- Examples include:
 - k-means clustering
 - DBSCAN
 - LOF (Local Outlier Factor)
 - Isolation Forest

Challenges:

- 1. Assumption might not always hold:
 - What if normal data is diverse and doesn't cluster well?

• Some applications (e.g., intrusion detection) have very varied normal activity.

In this case:

- Clustering methods may misclassify **normal but diverse** instances as outliers.
- Outliers may form small, tight groups that resemble normal clusters → Undetected.

2. High false positive rate:

 Normal points might appear "far" from others due to data spread, not because they're outliers.

***** Methods based on clustering:

- Detect clusters.
- Anything **not part of a cluster** = outlier.

But there are issues:

- A non-clustered point might just be **noise**, not an actual outlier.
- Cluster detection is computationally expensive.

Latest improvements:

New methods try to:

- Skip full clustering.
- Use proximity, density, or isolation-based metrics to directly detect outliers.

These are covered in later sections like 12.4 (Proximity-based) and 12.5 (Clustering-based).

✓ Summary:

- Strength: No need for labeled data.
- Weakness: May struggle when normal data lacks structure.

Real-life examples:

- Intrusion detection in networks
- Detecting faults in industrial sensors

✓ Semi-Supervised Methods

Definition: A middle ground between supervised and unsupervised methods.

- You have a **small number of labeled examples** (often only normals or only outliers).
- Most of the data is unlabeled.

★ Common cases:

- 1. Labeled normal objects only:
 - Build a model of normal behavior.
 - Any data point that doesn't fit → Outlier.
- 2. Labeled outliers only:
 - Harder. A few labeled outliers might not represent all outlier types.
 - Risk: Incomplete model of abnormality.

Strategy:

- Combine semi-supervised learning (e.g., self-training, graph-based methods) with outlier detection.
- Learn models from labeled normals, or use **unsupervised clustering** to help when few labels are available.

Challenges:

- Hard to model outliers when you only have a few labeled ones.
- Use unsupervised methods to help model normal data.
- Combine both sources for stronger detection.

✓ Summary:

- Strength: Leverages limited labeled data effectively.
- Weakness: Hard to generalize from small outlier samples.

Real-life examples:

- Medical diagnosis using partial diagnosis data
- System monitoring with known good behavior but rare fault examples



Overall Comparison Table

Method	Labeled Data Needed	Assumptions Made	Strengths	Weaknesses
Supervised	Yes – labeled normals/outliers	Model normal and/or abnormal behavior	High accuracy with good labels	Hard to get labeled outliers; class imbalance
Unsupervised	No	Normal data is clustered or patterned	Works without labels	Can fail if normal data is diverse
Semi- Supervised	Some labeled data	Partial supervision helps detect outliers	Leverages both labeled & unlabeled data	Hard to model outliers with few examples



Final Notes

- The choice of outlier detection method depends on the data available and the structure of normal and outlier behaviors.
- No method works best for all cases. Understanding the data, the application context, and available labels is critical for success.

Would you like me to summarize the next part: 12.3 Statistical Approaches or explain examples of each method in real-world use cases like fraud detection, manufacturing faults, or cybersecurity?

12.2.2 Statistical, Proximity-Based, and Clustering-Based Methods - Full **Explanation**

As discussed earlier, outlier detection depends on assumptions about the nature of normal data vs outliers. Based on the assumptions about data distribution and relationships, outlier detection techniques are generally divided into three categories:

1. Statistical Methods (Model-Based)

Key Assumption:

Normal data follows a known statistical distribution (e.g., Gaussian/normal distribution). Any data point that does not conform to the model is considered an outlier.

Core Idea:

If we assume the majority of data is generated from a probability distribution (e.g., Gaussian), we can:

- Estimate the **probability** that a data point belongs to this distribution.
- Points with very low probability are considered outliers.

***** Example 12.5 – Using a Gaussian Model:

In a 2D space (Figure 12.1), assume the data (excluding region R) fits a Gaussian distribution gD.

- For any data point x, gD(x) gives the **probability density** at x.
- For a point y in region R, the value gD(y) is very small \rightarrow meaning y likely does not **belong** to this distribution \rightarrow So, it is an **outlier**.

Model Types:

Parametric: Assumes a specific distribution form (e.g., normal, exponential).

- Needs to estimate parameters (mean, variance, etc.).
- Nonparametric: Makes fewer assumptions, often uses data-driven estimation.

Limitations:

- Assumes the data **follows a known distribution** → Not always true in real-world cases.
- If data deviates from assumed model, method may fail.

✓ Summary:

- Good for: Well-structured, statistically distributed data.
- Weakness: Not robust to distributional deviations.

2. Proximity-Based Methods

Key Assumption:

An object is an outlier if it is **far away** from its neighbors, **unlike most other points** in the data. These methods use **distance** or **density** to measure proximity.

Core Idea:

- Normal objects: Located close to their neighbors.
- Outliers: Have few or distant neighbors.

***** Example 12.6 – Using Proximity:

In Figure 12.1:

- Use each object's 3 nearest neighbors to measure proximity.
- Points in region R have 2nd and 3rd nearest neighbors that are much farther away than other points.
- Hence, those in R are proximity outliers.

Variants:

- **Distance-based**: Look at **distance** to k-nearest neighbors.
- **Density-based**: Compare the **local density** of a point to its neighbors (e.g., LOF Local Outlier Factor).

Limitations:

- 1. Choice of distance metric matters (Euclidean, Manhattan, etc.).
- 2. Difficult to define proximity for complex or high-dimensional data.
- 3. **Group outliers**: If outliers are **close to each other**, proximity-based methods may miss them (they don't look "different enough").

Summary:

- Good for: Spatial or numerical data where distance makes sense.
- Weakness: Fails if proximity is ill-defined or group outliers exist.

3. Clustering-Based Methods

Key Assumption:

- Normal objects belong to large, dense clusters.
- Outliers are:
 - In small/sparse clusters, or
 - Not part of any cluster.

Core Idea:

If most data forms clear clusters, then:

• Objects outside clusters, or in tiny clusters, are likely anomalies.

***** Example 12.7 – Using Clustering:

In Figure 12.1:

- Two clusters:
 - c1 (normal): Contains most points.
 - c2 (tiny): Just 2 points in region R.
- Since c2 is much smaller, it likely represents outliers.

*** How It Works:**

- 1. Run a clustering algorithm (e.g., k-means, DBSCAN).
- 2. Identify clusters and their sizes.
- 3. Points in:
 - Tiny/sparse clusters → Marked as outliers.
 - No cluster → Also marked as outliers.

Limitations:

- Clustering is computationally expensive, especially on large datasets.
- Outlier detection requires analyzing all data points, not just clusters.
- Might miss global outliers that appear normal within local clusters.

Summary:

- Good for: Data where clusters clearly separate normals from outliers.
- Weakness: Not scalable; depends on clustering success.

Final Comparison Table

Method Type	Assumption	Key Idea	Strengths	Weaknesses
Statistical	Data follows known distribution (e.g., Gaussian)	Model normality; flag low-probability points as outliers	Strong theoretical foundation	Not robust to wrong model assumptions
Proximity- Based	Outliers are far from their neighbors	Measure distance/density from neighbors	No label needed; intuitive	Sensitive to proximity definition; misses grou outliers
Clustering- Based	Normals form large clusters; outliers don't	Cluster data → flag small/sparse/out-of- cluster points	Leverages structure in data	Computationally expensive; sensitive to clustering quality

✓ Summary of Section 12.2.2

This section categorizes outlier detection methods based on the **assumptions about data behavior**, rather than data labeling. Each method has its **strengths and weaknesses**, and the best choice depends on:

- The data type (e.g., numerical, structured, distributed),
- The underlying distribution,
- The presence or absence of structure (clusters, patterns),
- And the scale of the dataset.

Would you like me to now explain **Section 12.3**: **Statistical Approaches in Detail** or go deeper into **distance-based vs density-based** proximity methods?