

# Final Notes & Review: Gene Expression Analysis & Clustering

## Part 1: The "Big Picture" - From Sequence to Action

This topic represents a major shift in focus. So far, we've treated the genome as a static "blueprint." Now, we're asking a dynamic question: **Which genes are "on" or "off" in a cell at a specific time, and how does this pattern of activity relate to the cell's function or condition (e.g., healthy vs. cancerous)?**

- **Jargon: Gene Expression**
  - **Simple Explanation:** This is the process of a gene being actively used to make a product, usually a protein. The *level* of gene expression is how much product is being made. A highly expressed gene is "loud," while a lowly expressed gene is "quiet."
- **The Goal:** We want to measure the expression levels of thousands of genes simultaneously. The resulting dataset is a massive table where **rows are genes** and **columns are different conditions or time points**. This is the **Gene Expression Matrix**.

## Part 2: The Technology - How We Measure Gene Expression

How do we listen to how "loud" each gene is?

- **The Proxy:** Measuring protein levels directly is difficult and complex. Instead, we measure the amount of **mRNA** for each gene. Since mRNA is the temporary "photocopy" of a gene used to make protein, its abundance is a good proxy for the gene's activity level.
- **The Two Main Techniques:**
  1. **Expression Microarrays (Gene Chips):**
    - **How it Works:** A glass slide is covered in millions of tiny spots. Each spot contains a known, unique DNA sequence called a **probe**, which corresponds to a specific gene. You extract all the mRNA from your cells, convert it to DNA (**cDNA**), and label it with a fluorescent dye. You then wash this glowing cDNA over the chip. The cDNA will **hybridize** (stick) to its complementary probe.
    - **The Result:** A spot corresponding to a highly expressed gene will have a lot of cDNA stuck to it and will glow very brightly. A spot for a lowly expressed gene will be dim. A camera measures the brightness of every spot to create the expression data.
  2. **RNA-Seq (The Modern Standard):**
    - **How it Works:** This is a more direct approach. You extract all the mRNA from your cells, convert it to cDNA, and then use the **Next-Generation Sequencing** technology (from the previous lecture) to sequence all of it.
    - **The Result:** You get millions of reads. You then map these reads back to the reference genome. The expression level of a gene is simply the **number of reads that map to that gene**. A highly expressed gene will produce many reads; a lowly expressed gene will produce few.

## Part 3: The Computational Challenge - Finding Patterns in the Data

We now have a massive matrix of numbers ( $n$  genes  $\times$   $m$  conditions). How do we make sense of it? This is where **Machine Learning**, specifically **Clustering**, comes in.

- **The Core Idea:** We are looking for hidden structures. The central hypothesis is that **genes that show similar expression patterns across various conditions are likely functionally related or co-regulated.**
  1. **Analogy:** Imagine tracking the stock prices of thousands of companies (genes) over a year (conditions). You would notice that all the airline stocks tend to go up and down together, and all the oil company stocks go up and down together. By clustering them based on their price patterns, you would have discovered the "airline sector" and the "oil sector" without anyone telling you they existed. We want to do the same for genes.
- **What Makes a Good Clustering? (Slide 11)**
  1. **Homogeneity:** Points within the same cluster should be "close" to each other.
  2. **Separation:** Points in different clusters should be "far" from each other.

#### Part 4: The Algorithms - How We Cluster the Genes

This lecture introduces two major families of clustering algorithms.

This is an algorithm that partitions data into a pre-defined number ( $k$ ) of clusters.

- **The Goal:** To find  $k$  center points (centroids) that minimize the **squared error distortion**—the average squared distance of every data point to its closest center.
- **The Lloyd Algorithm (Iterative Refinement):**
  - **Initialization:** Arbitrarily place  $k$  points to serve as the initial cluster centers.
  - **Assignment Step (Centers to Clusters):** Assign each data point (each gene) to the cluster of its *closest* center. This partitions all the data.
  - **Update Step (Clusters to Centers):** Recalculate the position of each of the  $k$  centers by finding the **center of gravity** (the mean or average) of all the data points assigned to it.
  - **Repeat:** Repeat steps 2 and 3 until the cluster centers stop moving.
- **Limitations of k-Means:**
  - **Gets Stuck:** It can easily get stuck in a **local minimum**; the final result depends heavily on the random initial placement of centers. (The **k-Means++ Initializer** is a smarter way to choose starting points to help with this).
  - **Spherical Assumption:** It naturally finds spherical, evenly-sized clusters and fails on complex shapes.
  - **"Hard" Assignments:** It forces every point to belong to exactly one cluster, which is problematic for points that lie between two clusters.

This algorithm doesn't produce a single clustering, but a tree of nested clusters.

- **The Goal:** To build a hierarchy of clusters from the bottom up (agglomerative).
- **The Algorithm:**
  - **Initialization:** Start with every data point (gene) as its own tiny cluster.
  - **Distance Matrix:** Calculate the distance between every pair of points.

- **Iteration:**
    - Find the **two closest clusters** in the entire dataset.
    - **Merge** them into a new, single cluster.
    - Update the distance matrix to include this new cluster.
  - **Repeat:** Repeat step 3 until only one giant cluster (containing all the points) remains.
  - **The Output: A Dendrogram.** The result is a tree diagram called a dendrogram. It visually shows the merge history. You can then "cut" the dendrogram at a certain height to produce any number of clusters you want.
  - **Key Challenge:** How do you define the distance between two *clusters*?
    - **D\_min (Single Linkage):** The distance is the shortest distance between *any* two points in the two clusters.
    - **D\_avg (Average Linkage):** The distance is the average distance between *all* pairs of points in the two clusters.
- 

## Possible Exam Questions & Key Concepts

- **Q: Compare and contrast Microarrays and RNA-Seq for measuring gene expression.**
  - **A: Microarrays** are a hybridization-based method that measures the *relative* fluorescence of known probes. They are cheaper but can only measure known genes and are less sensitive. **RNA-Seq** is a sequencing-based method that *counts* the number of reads for each gene. It is more expensive but more accurate, can discover new genes/transcripts, and provides a direct digital count of expression.
- **Q: What is a gene expression vector?**
  - **A:** It is a single row in the gene expression matrix. It represents the expression levels of a single gene across all measured conditions or time points. We cluster these vectors to find co-regulated genes.
- **Q: What is "hard" vs. "soft" clustering distinction?**
  - **A: k-Means** is a **hard clustering** algorithm because it assigns each data point to exactly one cluster. The **EM Algorithm** (for soft k-Means) is a **soft clustering** algorithm because it assigns each data point a "responsibility" or probability of belonging to *every* cluster.
- **Q: You run the k-Means algorithm twice on the same dataset and get two different results. Why did this happen?**
  - **A:** This happened because k-Means is sensitive to its **initialization**. The algorithm starts by randomly placing the initial k centers. Different random starting points can cause the algorithm to converge to different **local minima**, resulting in different final clusterings.
- **Q: What is a dendrogram and which clustering algorithm produces it?**
  - **A:** A dendrogram is a tree diagram that shows the nested hierarchy of clusters. It is produced by **Hierarchical Clustering**. The user can "cut" the dendrogram at a desired level to obtain a specific number of clusters.
- **Q: Explain the two key properties of a "good" clustering.**
  - **A:** 1) **Homogeneity:** Points within the same cluster should be close together and highly similar. 2) **Separation:** Points in different clusters should be far apart and dissimilar.

Of course. You are right to focus on the details, as that's where exam questions often test true understanding. Let's add that final, crucial layer of detail, covering the math, formulas, edge cases, and a direct comparison table.

---

## Final Deep Dive: Math, Formulas, and Edge Cases

### Part 1: The Math of k-Means Clustering

This section translates the abstract goal of "good clustering" into a precise mathematical optimization problem.

#### 1. The Distance Function: $d(v, C)$ (Distance from a Point to Centers)

- **Formula:**  $d(v, C) = \min( d(v, C_i) )$  for all centers  $C_i$  in the set  $C$ .
- **What it means:** The distance from a single data point  $v$  (a gene) to the entire set of cluster centers  $C$  is simply the Euclidean distance to its **single closest center**. This is the core of the "assignment" step.

#### 2. The Objective Function: $\text{distortion}(\text{Data}, C)$

- **What it is:** This is the function we are trying to **minimize**. It gives us a single number that tells us how "bad" our current choice of cluster centers is. A lower value is better.
- **Formula:**  $\text{distortion}(\text{Data}, C) = (1/n) * \sum [d(p, C)]^2$  for all  $n$  points  $p$  in the Data.
- **In Plain English:**
  1. For every single data point ( $p$ ), find its distance to its nearest cluster center ( $d(p, C)$ ).
  2. Square that distance.
  3. Sum up all these squared distances for all the data points.
  4. (Optional but common) Divide by the total number of points ( $n$ ) to get the mean squared error.
- **The Goal of k-Means:** Find the set of  $k$  center points  $C$  that makes this distortion value as small as possible.

### Part 2: The Math of Hierarchical Clustering

This algorithm doesn't minimize a single global function. Instead, it makes a series of locally optimal decisions based on a distance metric between clusters.

**The Key Question:** How do you define the distance  $D(C_1, C_2)$  between two *clusters* of points?

- **Formula 1:  $D_{\min}$  (Single-Linkage)**
  - $D_{\min}(C_1, C_2) = \min( d(i, j) )$  for all points  $i$  in  $C_1$  and  $j$  in  $C_2$ .

- **In Plain English:** The distance between two clusters is the distance between their **two closest points**.
  - **Effect:** This method is good at finding long, chained-out, or non-spherical shapes. However, it is very sensitive to outliers and can sometimes merge two otherwise distinct clusters just because a single point from each are close (this is called the "chaining" effect).
  - **Formula 2: D\_avg (Average-Linkage)**
    - $D\_avg(C\_1, C\_2) = (1 / |C\_1| * |C\_2|) * \sum d(i, j)$  for all points  $i$  in  $C\_1$  and  $j$  in  $C\_2$ .
    - **In Plain English:** The distance is the **average distance** between *every possible pair* of points between the two clusters.
    - **Effect:** This is more robust and less sensitive to outliers than single-linkage. It tends to produce more compact, globular clusters. It is often a good default choice.
  - *(Other common methods exist, like D\_max (complete-linkage) and Ward's method, which minimizes the increase in variance upon merging.)*
- 

### Part 3: Algorithm Comparison Table

Feature	k-Means Clustering (Lloyd's)	Hierarchical Clustering
<b>Input</b>	The data and a fixed number of clusters, <b>k</b> .	The data and a linkage criterion (e.g., min, avg).
<b>Output</b>	A single, flat partition of the data into <b>k</b> clusters.	A <b>dendrogram</b> (tree) showing a hierarchy of nested clusters.
<b>How it Works</b>	Iteratively moves centers to the mean of their assigned points.	Iteratively merges the two closest clusters.
<b>Cluster Shape</b>	Assumes clusters are <b>spherical</b> and roughly equal in size.	Can handle <b>arbitrary shapes</b> (especially with single-linkage).

<b>Time Complexity</b>	$O(n * k * i * m)$ , where $n$ =points, $k$ =clusters, $i$ =iterations, $m$ =dimensions. Generally considered <b>fast</b> for large datasets.	At least $O(n^2 * m)$ . It is <b>slow</b> for large datasets because of the need to compute and update the $n \times n$ distance matrix.
<b>Big Advantage</b>	<b>Speed and scalability.</b> It's much faster than hierarchical clustering on large datasets.	<b>Provides a full hierarchy.</b> Does not require $k$ to be specified in advance and gives a rich, visual output of the data's structure.
<b>Big Disadvantage</b>	<b>Must specify <math>k</math> in advance.</b> Highly sensitive to the random initialization and can get stuck in local optima.	<b>Computationally expensive.</b> Prohibitively slow for very large $n$ . The results can be sensitive to the choice of linkage criterion.
<b>Use Case</b>	Best for <b>large datasets</b> when you have a good idea of how many clusters to expect and the clusters are likely globular (e.g., segmenting customers).	Best for <b>smaller datasets</b> when the underlying structure is unknown or hierarchical in nature (e.g., analyzing evolutionary relationships between species).

#### Part 4: Edge Cases and Risky Questions for the Exam

- **Q: You run k-Means with  $k=3$  on a dataset that clearly contains two long, thin, parallel clusters. What is the likely outcome?**
  - **A:** The algorithm will likely fail to identify the true clusters. Because k-Means tries to find spherical centers and minimize distance to those centers, it will probably cut each of the long clusters in half, creating three meaningless, mixed clusters that do not reflect the true structure of the data.
- **Q: What is the main drawback of single-linkage hierarchical clustering?**
  - **A:** It is highly sensitive to noise and outliers. A single noisy data point that happens to lie between two otherwise well-separated clusters can act as a "bridge," causing the algorithm to merge them prematurely. This is known as the **chaining effect**.
- **Q: How do you choose the value of  $k$  for the k-Means algorithm?**

- **A:** There is no single perfect method. It often involves a combination of domain knowledge and empirical methods. One common technique is the "**Elbow Method**," where you run k-Means for many different values of k and plot the distortion (the objective function). The "elbow" of the resulting curve—the point where the rate of decrease sharply slows—is often a good estimate for the optimal k.
- **Q: Your gene expression data has 10,000 genes (points) and 50 conditions (dimensions). Which clustering algorithm would be more computationally feasible to run?**
  - **A: k-Means.** Hierarchical clustering would require computing and storing a massive 10,000 x 10,000 distance matrix, which would be extremely slow and memory-intensive. K-Means, while iterative, scales much better with the number of data points (n) and would be the only feasible choice.
- **Q: What are the two main steps of the Lloyd's algorithm for k-Means?**
  - **A:** 1) The **Assignment Step**, where each data point is assigned to its closest current cluster center. 2) The **Update Step**, where each cluster center is moved to the mean (center of gravity) of all the points assigned to it. These two steps are repeated until convergence.

## Solutions & Explanations: Clustering Gene Expression Data

### Question 1: K-Means and Soft K-Means Simulation

**Data:**  $p_1=(1,2)$ ,  $p_2=(5,6)$ ,  $p_3=(4,3)$ ,  $p_4=(6,4)$

**Initialization:** Let's choose  $p_1$  and  $p_4$  as initial centers.  $C_1 = (1,2)$ ,  $C_2 = (6,4)$ .

#### A) K-Means Clustering (1 Iteration)

**1. Assignment Step:** Assign each point to its closest center. We need to calculate squared Euclidean distances:  $d^2 = (x_2 - x_1)^2 + (y_2 - y_1)^2$ .

- **Point  $p_1=(1,2)$ :**  $d^2(p_1, C_1) = 0$ .  $d^2(p_1, C_2) = (6-1)^2 + (4-2)^2 = 25+4 = 29$ . Closest to **C1**.
- **Point  $p_2=(5,6)$ :**  $d^2(p_2, C_1) = (1-5)^2 + (2-6)^2 = 16+16 = 32$ .  $d^2(p_2, C_2) = (6-5)^2 + (4-6)^2 = 1+4 = 5$ . Closest to **C2**.
- **Point  $p_3=(4,3)$ :**  $d^2(p_3, C_1) = (1-4)^2 + (2-3)^2 = 9+1 = 10$ .  $d^2(p_3, C_2) = (6-4)^2 + (4-3)^2 = 4+1 = 5$ . Closest to **C2**.

- **Point p4=(6,4):**  $d^2(p4, C1) = (1-6)^2 + (2-4)^2 = 25 + 4 = 29$ .  $d^2(p4, C2) = 0$ . Closest to **C2**.
- **Resulting Clusters:**
  - Cluster 1: {p1}
  - Cluster 2: {p2, p3, p4}

**2. Update Step:** Recalculate centers by finding the mean of the points in each cluster.

- **New C1:**  $\text{mean}(\{p1\}) = (1, 2)$ .
- **New C2:**  $\text{mean}(\{(5,6), (4,3), (6,4)\}) = ((5+4+6)/3, (6+3+4)/3) = (15/3, 13/3) = (5, 4.33)$ .

**End of Iteration 1:** The new cluster centers are  $C1=(1,2)$  and  $C2=(5, 4.33)$ .

## B) Soft K-Means Clustering (1 Iteration, E-step only)

We need to calculate the "responsibility" of each center for each point. The formula uses  $e^{(-\beta*d)}$ , where  $d$  is the Euclidean distance (not squared). Let's use the first formula from the slides which uses  $1/d^2$  as it is simpler to calculate from our work above.

### 1. E-Step: Calculate Responsibilities.

Responsibility of  $C_i$  for point  $p$ :  $R(C_i, p) = (1/d^2(p, C_i)) / \sum (1/d^2(p, C_j))$

- **For p2=(5,6):**  $d^2(p2, C1)=32$ ,  $d^2(p2, C2)=5$ .
  - Responsibility of C1:  $(1/32) / (1/32 + 1/5) = 0.03125 / (0.03125 + 0.2) = 0.135$
  - Responsibility of C2:  $(1/5) / (1/32 + 1/5) = 0.2 / 0.23125 = 0.865$
  - *(Checks out:  $0.135 + 0.865 = 1.0$ )*
- **For p3=(4,3):**  $d^2(p3, C1)=10$ ,  $d^2(p3, C2)=5$ .
  - Responsibility of C1:  $(1/10) / (1/10 + 1/5) = 0.1 / (0.1 + 0.2) = 0.333$
  - Responsibility of C2:  $(1/5) / (1/10 + 1/5) = 0.2 / 0.3 = 0.667$

**End of E-Step:** Soft K-Means doesn't assign points to a single cluster. Instead, it would say p2 belongs 86.5% to C2 and 13.5% to C1. p3 belongs 66.7% to C2 and 33.3% to C1. The M-step would then calculate new centers using these responsibilities as weights.

## Question 2: Probabilities for k-Means++ Initialization

If  $p1=(1,2)$  is the first center, we need to calculate the probability of picking each other point, which is proportional to its squared distance ( $D^2$ ) from the set of current centers (which is just  $p1$ ).

1. **Calculate  $D^2$  for remaining points:**
  - $D^2(p2) = d^2(p2, p1) = (5-1)^2 + (6-2)^2 = 16 + 16 = 32$
  - $D^2(p3) = d^2(p3, p1) = (4-1)^2 + (3-2)^2 = 9 + 1 = 10$
  - $D^2(p4) = d^2(p4, p1) = (6-1)^2 + (4-2)^2 = 25 + 4 = 29$
2. **Calculate Total  $D^2$  Sum:**
  - $\text{Sum} = 32 + 10 + 29 = 71$



### 3. Calculate Probabilities:

- $\text{Prob}(p_2) = D^2(p_2) / \text{Sum} = 32 / 71 \approx 45.1\%$
- $\text{Prob}(p_3) = D^2(p_3) / \text{Sum} = 10 / 71 \approx 14.1\%$
- $\text{Prob}(p_4) = D^2(p_4) / \text{Sum} = 29 / 71 \approx 40.8\%$

The k-Means++ initializer is more likely to pick points that are far away from existing centers, leading to a better initial placement.

---

### Question 3: Limitations of Hard K-Means and Mitigations

Limitation	Mitigation/Solution
<b>1. Sensitive to Initialization / Gets stuck in local optima.</b>	Use a smarter initialization algorithm like <b>k-Means++</b> , which chooses initial centers that are far apart. Also, run the algorithm multiple times with different random initializations and choose the result with the lowest distortion.
<b>2. Must specify the number of clusters, k, in advance.</b>	Use methods like the " <b>Elbow Method</b> " (plotting distortion vs. k) or use a different algorithm like <b>Hierarchical Clustering</b> which doesn't require k as an input.
<b>3. Assumes clusters are spherical and of similar size.</b>	If clusters have complex shapes (e.g., long and thin, or crescent-shaped), use a different algorithm like <b>DBSCAN</b> (density-based clustering) or <b>Hierarchical Clustering</b> with single linkage.
<b>4. Makes "hard" assignments, which is problematic for ambiguous points.</b>	Use a " <b>soft</b> " clustering algorithm like the <b>EM Algorithm (Soft K-Means)</b> , which assigns a probability of belonging to each cluster instead of a hard assignment.

---

## Question 4 & 5: Biological Motivations

### Why is clustering genes important in Biology?

Clustering is fundamental for discovering hidden relationships in complex gene expression data. Its primary importance lies in the "**guilt-by-association**" principle:

1. **Function Discovery:** Genes that show similar expression patterns (i.e., cluster together) are often involved in the same biological pathway or cellular process. If you have a cluster of 10 genes where 9 are known to be involved in metabolism, you can form a strong hypothesis that the 10th unknown gene is also a metabolism gene.
2. **Regulatory Discovery:** Genes that cluster together are often controlled by the same **transcription factors**. Finding these co-regulated gene sets is the first step toward identifying the regulatory motifs in their upstream regions.

### What are the biological motivations for Soft K-Means and Hierarchical Clustering?

- **Soft K-Means:** Biology is not always black and white. A single gene can sometimes play a role in **multiple biological pathways**. A hard clustering algorithm like k-Means forces you to assign this gene to just one group. Soft K-Means is biologically more realistic because it can assign a gene partial "responsibility" to multiple clusters, reflecting its potential involvement in several processes.
  - **Hierarchical Clustering:** Biological processes are often inherently hierarchical. For example, "metabolism" is a large category, which contains smaller sub-categories like "carbohydrate metabolism" and "lipid metabolism." Hierarchical clustering is perfectly suited to discover and represent this nested structure. The resulting dendrogram provides a rich, multi-level view of gene relationships that a flat clustering from k-Means cannot.
-

## Question 6 & 7: Hierarchical Clustering Simulation

### Data:

GeneA: [2.0, 3.0, 2.5, 3.5]

GeneB: [2.2, 3.1, 2.7, 3.6]

GeneC: [5.0, 5.5, 5.1, 5.4]

GeneD: [8.0, 7.5, 7.8, 7.7]

### Step 1: Calculate the initial distance matrix (Euclidean Distance).

- $d(A,B) = \sqrt{(2.2-2)^2+(3.1-3)^2+(2.7-2.5)^2+(3.6-3.5)^2} = \sqrt{0.04+0.01+0.04+0.01} = \sqrt{0.1} \approx 0.316$
- $d(A,C) = \sqrt{(5-2)^2+\dots} \approx \sqrt{9 + 6.25 + 6.76 + 3.61} \approx 5.06$
- ...and so on. The matrix looks like:

	A	B	C	D
A	0	<b>0.316</b>	5.06	8.87
B		0	4.88	8.64
C			0	4.14
D				0

### A) Single Linkage Method

- **Iteration 1:** The smallest distance is **0.316** between A and B. We merge them into a new cluster (AB).
- **Iteration 2:** We update the distance matrix. The distance from a point X to (AB) is  $\min(d(X,A), d(X,B))$ .
  - $d(C, (AB)) = \min(d(C,A), d(C,B)) = \min(5.06, 4.88) = 4.88$
  - $d(D, (AB)) = \min(d(D,A), d(D,B)) = \min(8.87, 8.64) = 8.64$
  - The only other distance is  $d(C,D) = 4.14$ .
  - The smallest distance is now **4.14** between C and D. We merge them into (CD).
- **Iteration 3:** We have two clusters left: (AB) and (CD). The distance is  $\min(d(A,C), d(A,D), d(B,C), d(B,D)) = \min(5.06, 8.87, 4.88, 8.64) = 4.88$ . We merge them.
- **Final Dendrogram (Single Linkage):** ((A,B), (C,D)) with A and B merging first.

### B) Average Linkage Method

- **Iteration 1:** Same as single linkage. We merge A and B into (AB).
- **Iteration 2:** We update the distance matrix. The distance from X to (AB) is  $\text{avg}(d(X,A), d(X,B))$ .
  - $d(C, (AB)) = (d(C,A) + d(C,B)) / 2 = (5.06 + 4.88) / 2 = 4.97$
  - $d(D, (AB)) = (d(D,A) + d(D,B)) / 2 = (8.87 + 8.64) / 2 = 8.755$
  - $d(C,D) = 4.14$ .
  - The smallest distance is still **4.14** between C and D. We merge them into (CD).
- **Iteration 3:** The distance between (AB) and (CD) is the average of all four cross-cluster distances:  $(d(A,C)+d(A,D)+d(B,C)+d(B,D))/4 = (5.06+8.87+4.88+8.64)/4 = 6.86$ .

- **Final Dendrogram (Average Linkage):** The structure is the same, ((A,B), (C,D)), but the merge distances recorded in the dendrogram would be different.
- 

**Question 8: What are the possible disadvantages of hierarchical clustering?**

1. **High Computational Cost:** It has a time complexity of at least  $O(n^2)$ , making it very slow and memory-intensive for datasets with a large number of points (genes).
2. **Irreversibility (Greedy Nature):** Once a decision to merge two clusters is made, it can never be undone. A single bad merge early on (e.g., due to an outlier) can lead to a poor overall clustering structure.
3. **Ambiguity in Defining Cluster Distance:** The final result can look very different depending on the choice of linkage method (single, average, complete, etc.), and it's not always obvious which method is best for a given dataset.
4. **Dendrogram Interpretation:** Interpreting the dendrogram can be subjective. Deciding where to "cut" the tree to define the final clusters is often arbitrary and can lead to different conclusions.

–Fahad Nadim Ziad, 24341216