

Lesson 30 Summary: Hierarchical Clustering

Data Science with Python – Key Concepts

Data Science Program

Hierarchical Clustering

Agglomerative

Bottom-up approach
Merge closest pairs
Most common type

Linkage Methods

Single: nearest points
Complete: farthest
Ward: minimize variance

Dendrogram

Tree visualization
Cut at height for K
Shows hierarchy

Advantages

No need to pre-specify K
Interpretable structure

Limitations

$O(n^2)$ memory | Slow for large data
No reassignment after merge

`dendrogram(linkage(X, method="ward"))` | `AgglomerativeClustering(n_clusters=K)`

Hierarchical clustering builds a tree of clusters

Bottom-up clustering:

- **Start:** Each point is its own cluster
- **Merge:** Combine two closest clusters
- **Repeat:** Until single cluster or threshold

Agglomerative is the most common hierarchical method

How to measure cluster distance:

- **Single:** Minimum distance between points
- **Complete:** Maximum distance between points
- **Average:** Mean distance between all pairs
- **Ward:** Minimizes within-cluster variance

Ward linkage usually gives the most balanced clusters

The Dendrogram

Tree visualization:

- **X-axis:** Data points or cluster labels
- **Y-axis:** Distance at which clusters merge
- **Cut:** Horizontal line determines K clusters

Dendrograms reveal hierarchical structure

Choosing Number of Clusters

Methods:

- **Dendrogram:** Cut where gaps are large
- **Silhouette:** Same as K-means
- **Domain knowledge:** Business requirements

Dendrograms let you explore multiple K values visually

When to use hierarchical:

- **No pre-specified K:** Explore structure first
- **Interpretable:** Tree shows relationships
- **Small data:** Works well for hundreds of points

Great for exploratory analysis

Computational challenges:

- **Memory:** $O(n^2)$ distance matrix
- **Speed:** Slower than K-means
- **Irreversible:** Once merged, cannot undo

Not suitable for very large datasets

Two approaches:

- **scipy**: For dendrogram visualization
- **sklearn**: For prediction and labels

scipy for visualization, **sklearn** for clustering

Use `scipy.cluster.hierarchy` for dendrograms

Essential Commands:

Task	Code
Linkage matrix	<code>linkage(X, method='ward')</code>
Dendrogram	<code>dendrogram(Z)</code>
Get labels	<code>fcluster(Z, t=K, criterion='maxclust')</code>
sklearn	<code>AgglomerativeClustering(n_clusters=K)</code>

Hierarchical clustering reveals data structure