

Clustering Algorithm Complexity & Performance Guide

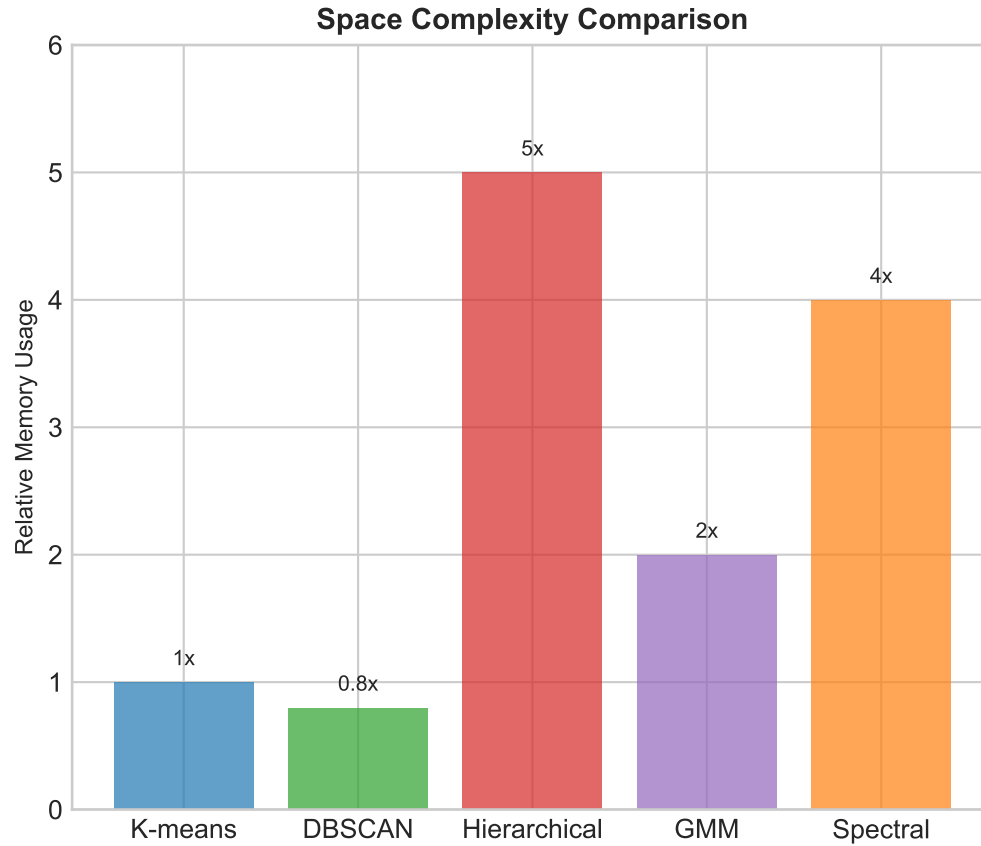
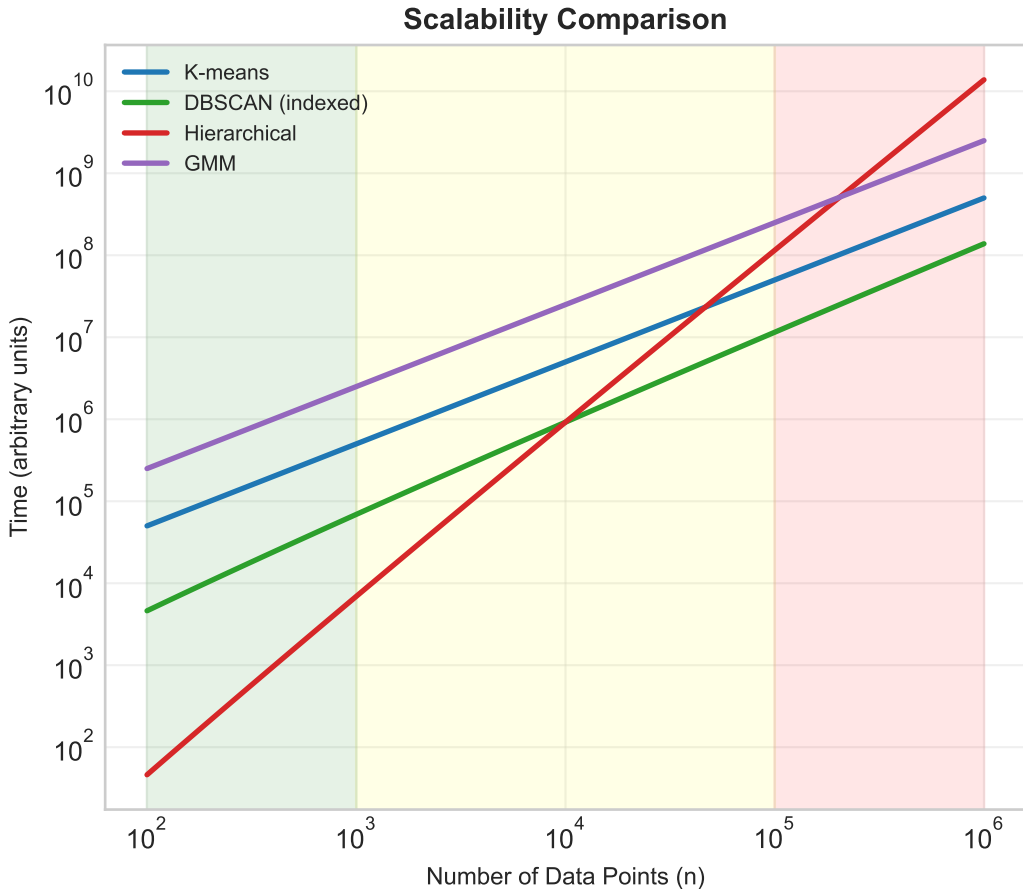
Algorithm Complexity Analysis

Big O Notation Comparison

Algorithm	Time Complexity	Space Complexity	Scalability
K-means	$O(n \cdot k \cdot i \cdot d)$	$O(n \cdot d + k \cdot d)$	Excellent
DBSCAN	$O(n^2)$ / $O(n \log n)^*$	$O(n)$	Good
Hierarchical	$O(n^3)$ / $O(n^2 \log n)^*$	$O(n^2)$	Poor
GMM	$O(n \cdot k^2 \cdot i \cdot d)$	$O(k \cdot d^2)$	Moderate
	$O(n^3)$	$O(n^2)$	Poor

Notation Guide:

n = number of data points
k = number of clusters
i = number of iterations
d = number of dimensions
* = with spatial index



Practical Recommendations

Small Data (<10K points)

→ Any algorithm works

Medium Data (10K-100K)

→ K-means or DBSCAN

Large Data (>100K)

→ MiniBatch K-means

High Dimensions (>50)

→ Consider PCA first

Real-time Requirements

→ Pre-computed K-means

Memory Constrained

→ Avoid Hierarchical

Optimization Techniques

MiniBatch K-means:

- Samples subset of data
- 10-100x faster on large data

Spatial Indexing (DBSCAN):

- KD-tree or Ball-tree
- $O(n^2) \rightarrow O(n \log n)$

Dimensionality Reduction:

- PCA before clustering
- Reduces d in $O(n \cdot k \cdot i \cdot d)$

Early Stopping:

- Monitor convergence
- Stop when stable

Implementation Complexity

Algorithm	Ease	Lines of Code*	Tuning
K-means	Easy	~50	Simple
DBSCAN	Moderate	~100	Tricky
Hierarchical	Easy	~30	Simple
GMM	Hard	~200	Complex
Spectral	Hard	~150	Complex

Performance Tips:

1. Profile before optimizing
2. Use vectorized operations
3. Consider approximate methods
4. Parallelize when possible

Choose algorithms based on data size, dimensionality, and performance requirements