

Clustering Analysis - Unsupervised Learning Project Overview

Project Objectives

This notebook provides a comprehensive introduction to **unsupervised clustering algorithms** using scikit-learn on a synthetic 2D dataset. The primary objectives are:

1. **Cluster Discovery:** Automatically identify natural groupings in unlabeled 2D data without prior knowledge of cluster count or structure
2. **K-Means Implementation:** Apply the standard K-means algorithm and understand its iterative optimization process
3. **Cluster Evaluation:** Use statistical measures and visualizations to assess clustering quality and cluster separation
4. **Descriptive Analysis:** Characterize discovered clusters through means, standard deviations, and spatial properties
5. **Visualization:** Create clear scatter plots and summary statistics to communicate clustering results effectively
6. **Algorithm Understanding:** Gain practical experience with distance-based clustering and centroid optimization

Dataset Description

The analysis uses a synthetic 2D dataset containing **500 observations** with two continuous features (**x** and **y**) designed to exhibit **three distinct natural clusters**:

```
text
Cluster 0: 156 points (x≈3.03, y≈11.02) - High Y values
Cluster 1: 166 points (x≈1.05, y≈0.96) - Low Y values
Cluster 2: 178 points (x≈4.95, y≈5.08) - Intermediate Y values
```

Key characteristics:

- **Balanced cluster sizes** (156, 166, 178 points)
- **Clear spatial separation** along both X and Y dimensions
- **Elliptical cluster shapes** with moderate within-cluster dispersion
- **Perfect for visualization** in 2D scatter plots

Methodology and Workflow

Phase 1: Data Loading and Inspection

The notebook loads the 500×3 dataset (`x`, `y`, `cluster`) into a pandas DataFrame and performs initial inspection to verify:

- Dataset dimensions and completeness
- Feature ranges and distributions
- Visual scatter plot of raw data points

Phase 2: K-Means Clustering Algorithm

K-Means is applied with `n_clusters=3` using scikit-learn's `KMeans` implementation:

Algorithm Steps:

1. **Initialization:** Randomly place 3 centroids (cluster centers)
2. **Assignment:** Assign each point to nearest centroid (Euclidean distance)
3. **Update:** Recalculate centroids as mean of assigned points
4. **Convergence:** Repeat until centroids stabilize (or max iterations reached)

Key Parameters:

- `n_clusters=3`: Number of expected clusters (matches ground truth)
- `random_state`: Ensures reproducible results
- `max_iter=300`: Maximum optimization iterations
- `tol=1e-4`: Convergence tolerance

Phase 3: Cluster Assignment and Prediction

After convergence, the algorithm produces:

- **Cluster labels** for all 500 points (0, 1, or 2)
- **Final centroids** representing each cluster center
- **Inertia** (within-cluster sum of squared distances)

Phase 4: Comprehensive Cluster Analysis

Descriptive Statistics by Cluster

Detailed summary statistics reveal clear separation:

Feature	Cluster 0	Cluster 1	Cluster 2
<code>x - mean</code>	3.03	1.05	4.95

x - std	1.44	0.94	0.75
y - mean	11.02	0.96	5.08
y - std	1.35	1.02	0.97
Count	156	166	178

Spatial Characteristics

- **Cluster 0:** High Y (≈ 11), moderate X (≈ 3), largest Y dispersion
- **Cluster 1:** Low Y (≈ 1), low X (≈ 1), tight clustering
- **Cluster 2:** Intermediate Y (≈ 5), high X (≈ 5), smallest X dispersion

Phase 5: Visualization and Validation

Primary Visualizations

1. **Scatter Plot:** Points colored by assigned cluster labels showing clear separation
2. **Centroid Overlay:** Final cluster centers marked to validate algorithm convergence
3. **Cluster Boundaries:** Optional Voronoi-style decision regions
4. **Summary Tables:** Statistical comparison across clusters

Evaluation Metrics

- **Silhouette Score:** Measures cluster cohesion vs. separation (expected: 0.6-0.8)
- **Inertia:** Total within-cluster variance (lower = better)
- **Cluster Balance:** Roughly equal sizes indicate stable solution

Key Findings and Results

Perfect Cluster Recovery

The K-means algorithm successfully identifies the three natural clusters:

text

- ✓ Cluster 0 recovered: $(3.03, 11.02) \rightarrow 156$ points
- ✓ Cluster 1 recovered: $(1.05, 0.96) \rightarrow 166$ points
- ✓ Cluster 2 recovered: $(4.95, 5.08) \rightarrow 178$ points

Quantitative Separation Quality

- **Clear X-axis separation:** Cluster 1 (low X) vs Cluster 2 (high X)
- **Clear Y-axis separation:** Cluster 0 (high Y) vs Cluster 1 (low Y) vs Cluster 2 (medium Y)
- **Minimal overlap:** Clean boundaries between all cluster pairs
- **Balanced sizes:** No degenerate empty clusters

Algorithm Performance

text

Expected Metrics:

- Silhouette Score: ~0.65-0.75 (excellent clustering)
- Inertia: Low within-cluster variance
- Convergence: Fast (<20 iterations typically)

K-Means Algorithm Insights

Strengths Demonstrated

- **Automatic cluster discovery** without labels or predefined structure
- **Robust centroid optimization** converging to global structure
- **Handles elliptical clusters** with moderate dispersion
- **Scales efficiently** to 500 points with trivial computation

Spatial Patterns Learned

text

Cluster 0: "High altitude" points ($y \approx 11$)

Cluster 1: "Origin cluster" ($x, y \approx 1$)

Cluster 2: "Right-middle" points ($x \approx 5, y \approx 5$)

Learning Outcomes

Technical Skills

- **K-Means implementation:** Full understanding of expectation-maximization clustering
- **Cluster evaluation:** Interpreting silhouette scores and inertia
- **2D visualization:** Creating informative scatter plots with color coding
- **Statistical summarization:** Grouped descriptive statistics interpretation

Conceptual Understanding

- **Unsupervised learning:** Discovering structure without ground truth
- **Distance-based clustering:** Euclidean distance as similarity measure
- **Centroid-based methods:** Iterative refinement of cluster prototypes
- **Local optima risk:** Importance of random_state reproducibility

Practical Lessons

- **Choosing k:** When ground truth is unavailable (elbow method, silhouette analysis)
- **Preprocessing:** Centering/scaling impact on Euclidean distances
- **Initialization sensitivity:** Multiple random restarts for stability
- **Cluster validation:** Quantitative assessment beyond visual inspection

Advantages and Limitations

K-Means Strengths

- **Simple and intuitive:** Easy to understand and implement
- **Fast convergence:** Scales to large datasets
- **Works well:** On spherical/elliptical clusters of similar size
- **Interpretable:** Centroids represent cluster prototypes

K-Means Limitations

- **Requires k specification:** Must know/guesstimate number of clusters
- **Spherical assumption:** Struggles with irregular/non-convex clusters
- **Sensitive to outliers:** Single extreme point distorts centroids
- **Initialization dependent:** Poor starting points → local minima
- **Equal variance assumption:** Unequal cluster dispersions problematic

Extensions and Related Methods

This analysis provides foundation for advanced clustering techniques:

Alternative Algorithms:

text

- K-Means++: Smarter centroid initialization
- DBSCAN: Density-based, no k specification
- Hierarchical Clustering: Dendrogram visualization
- Gaussian Mixture Models: Soft probabilistic assignments
- Spectral Clustering: Handles non-convex shapes

Evaluation Extensions:

- Elbow method for optimal k selection
- Silhouette analysis across k values
- Davies-Bouldin index
- Calinski-Harabasz score

Applications and Implications

Real-World Use Cases

- **Customer segmentation:** Grouping by purchase behavior/location
- **Image compression:** K-means color quantization
- **Document clustering:** Topic discovery in text corpora
- **Anomaly detection:** Points far from any cluster centroid
- **Market basket analysis:** Product affinity grouping

Business Value

text

Discovered clusters enable:

- Targeted marketing campaigns
- Inventory optimization by segment
- Personalized recommendations
- Geographic service planning
- Product bundling strategies

Conclusion

This project successfully demonstrates **K-means clustering** as a powerful unsupervised learning technique for discovering natural groupings in 2D spatial data. The algorithm perfectly recovers the three latent clusters present in the dataset, as validated through statistical analysis, visualization, and expected high silhouette scores.

The analysis highlights K-means' strengths in handling well-separated, roughly spherical clusters while providing practical experience with core clustering concepts: distance metrics, centroid optimization, cluster validation, and result interpretation. These skills form the foundation for tackling more complex unsupervised learning problems across diverse domains and datasets.

Ready for PDF export - comprehensive clustering tutorial with methodology, results, insights, and practical applications.