

## UL task: SpectralClustering

**Motivation:** Group Olivetti face images into identity-consistent clusters **without labels**. The dataset is high-dimensional and lies on a curved manifold (pose, lighting, expression). Clustering exposes this intrinsic structure directly from features.

### Datasets (design choice & rationale):

- **Olivetti Faces (real-world, images):** 400 images ( $64 \times 64 \rightarrow 4096\text{-D}$ ). High-D, small-n; benefits from **PCA** and is a realistic testbed for face embeddings.

### Model fit:

- **HDBSCAN:** after PCA ( $\approx 100\text{--}200$ ), can find identity pockets and mark outliers (`min_cluster_size` tunes strictness).
- **OPTICS:** density ordering shows structure; reasonable on PCA features.
- **Spectral (ours):** with `affinity='nearest_neighbors'` often recovers  $\approx 40$  subjects on PCA features; captures **non-linear** separations.

(We initially considered the generated Two-Moons dataset to illustrate non-convex structure, but we ultimately report only Olivetti for a purely real-world evaluation.)

### Metrics:

1. **Silhouette Coefficient (unsupervised)** — average “closer to own cluster than the next best” score using distances in feature space;  $[-1,1][\text{--}1,1]$ , higher is better.

**Why here:** lets us tune Spectral’s `n_neighbors` and `n_clusters` **without labels** on the PCA space. (Caveat: sensitive to cluster-size imbalance and large  $kk$ ; we report overall values.)

2. **Adjusted Rand Index (ARI) (external)** — chance-adjusted agreement with ground-truth identities;  $[-1,1][\text{--}1,1]$  with  $0 \approx 0 \approx$  random and  $1 = 1 =$  perfect.

**Why here:** Olivetti ships labels; ARI quantifies identity consistency post-hoc and enables a fair comparison to density methods. (Caveat: ARI tends to prefer matching the true number of clusters; we sweep  $kk$  around  $\approx 40$  and report both ARI and Silhouette.)

CI.3

### Model: Spectral Clustering

**Motivation:** Spectral converts clustering into **graph partitioning**. It builds a **k-nearest-neighbors (kNN)** similarity graph and clusters in the space of the **graph Laplacian's** leading eigenvectors, which captures **non-linear** structure (pose/lighting/expression) better than centroid-based models.

#### Key characteristics:

1. Handles **non-linear** cluster geometry via a neighborhood graph.
2. Less biased toward spherical clusters than k-means/GMM; tolerates shape/size imbalance.
3. Requires `n_clusters` —useful here because prior knowledge is  $\approx 40$  subjects.

CI.4

### How Spectral Clustering works (method, optimization, implementation)

**Step 1 – Affinity graph:** Build sparse **kNN graph**  $WW$  (edges between nearest neighbors).

**Step 2 – Normalized Laplacian:**  $L_{sym} = I - D^{-1/2} W D^{-1/2}$ ,  $L_{sym} = I - D^{-1/2} W D^{-1/2}$ , with  $D_{ii} = \sum_j W_{ij}$ ,  $D_{ii} = \sum_j W_{ij}$ .

**Step 3 – Spectral embedding:** Take the  $kk$  eigenvectors of  $L_{sym}L_{sym}$  with **smallest** eigenvalues; row-normalize to get  $U \in \mathbb{R}^{n \times k}$ ,  $U \in \mathbb{R}^{n \times k}$ .

**Step 4 – Discretization:** Run **k-means** on rows of  $U$  to obtain  $kk$  clusters.

**Optimization view:** Spectral is the relaxed solution of **Normalized Cuts**

$$Ncut(A_1, \dots, A_k) = \sum_{r=1}^k kcut(A_r, A^T r) \text{assoc}(A_r, V).$$

$$Ncut(A_1, \dots, A_k) = \sum_{r=1}^k k \text{assoc}(A_r, V) \text{cut}(A_r, A^T r).$$

The exact minimization is NP-hard; the **spectral relaxation** reduces it to an eigenproblem on the Laplacian, then k-means discretizes the relaxed solution.

#### Our implementation (scikit-learn):

- Preprocessing: **PCA(100, whiten=True)** on flattened images (denoise + stable distances).
- Model: `SpectralClustering(affinity="nearest_neighbors", assign_labels="kmeans")`.

- Swept hyper-params: `n_clusters ∈ {30,35,38,40,42,45,50}` , `n_neighbors ∈ {8,10,12,15,20}` , `random_state=42` .
- 

CI.5

### Results for Spectral Clustering (Olivetti Faces)

**Input:**  $n=400$ , 40 ground-truth subjects; features used: **PCA(100)**  $\rightarrow X \in \mathbb{R}^{400 \times 100} \times \mathbb{R}^{400 \times 100}$

---

**Parameters:** `n_clusters=38` , `n_neighbors=10` **Evaluation**

**metrics:silhouette\_score:** **0.053** **adjusted\_rand\_score:** **0.272** (*best ARI in our sweep*) (see figure: *spectral\_olivetti\_pca100\_k38\_nn10.png*)

**Parameters:** `n_clusters=40` , `n_neighbors=12`

**Evaluation metrics:**

- **silhouette\_score:** **0.047**
- **adjusted\_rand\_score:** **0.199**

(see figure: *spectral\_olivetti\_pca100\_k40\_nn12.png*)

**Discussion & findings:**

- **Best setting by ARI:** ( $k=38$ ,  $nn=10$ ) slightly **under-specifies** the nominal 40 identities, avoiding over-splitting and yielding the best identity consistency (**ARI**  $\approx 0.272$ ).
  - **Silhouette magnitudes:** **low (~0.05)** across settings  $\Rightarrow$  clusters are **partly overlapping** in PCA(100) with Euclidean distance; identity separation exists but is weak in this space.
  - **Effect of `n_neighbors`:** Too small fragments the graph; too large over-smooths and merges identities. The sweet spot is **10–12 neighbors**.
  - **Full assignment:** Spectral labels **every** sample (no `1` noise), which simplifies Silhouette reporting and contrasts with density methods that may leave hard cases as noise.
- 

CI.6

**Summary & takeaways (Spectral on Olivetti):**

Spectral Clustering on **PCA(100)** features provides a principled non-linear partition with **moderate identity consistency**(best **ARI  $\approx 0.272$** , **Silhouette  $\approx 0.05$** ). The model is most effective with a **kNN affinity** ( $\approx 10\text{--}12$  neighbors) and  $k$  near—but not necessarily equal to—the subject count. The low Silhouette indicates that Euclidean geometry on PCA(100) remains **blurry** for identities; stronger embeddings (more PCs, non-linear DR, or learned features) are likely to improve separation. Overall, Spectral offers a clean, deterministic baseline for comparison against density-based methods on real-world face data.