

STAT 5244 – Unsupervised Learning

Homework 1

Name: Chuyang Su UNI: cs4570

1 Dimension Reduction on Digits Data.

1.1 Apply linear dimension reduction techniques.

In this experiment, I applied three linear dimension reduction methods and compared their performance on the `scikit-learn` *Digits* dataset ($n = 1797$, $p = 64$).

Each method projects the data into a two-dimensional latent space, on which I visualized the results and quantitatively evaluated their ability to separate the ten digit classes.

The results of this experiment are summarized below.

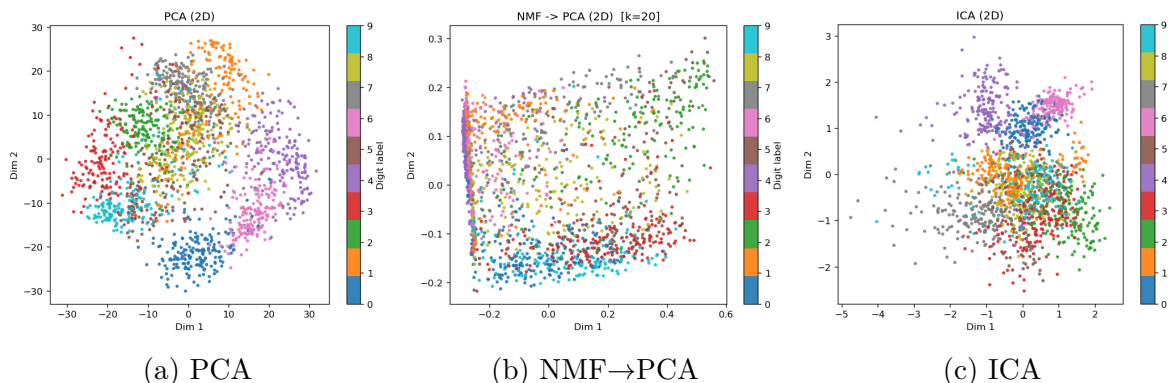


Figure 1: 2D embeddings of the digits data using PCA, NMF→PCA, and ICA. Colors denote true digit labels.

Method	ARI	NMI	Silhouette
PCA	0.3614	0.5190	0.3993
NMF → PCA	0.1588	0.2961	0.4080
ICA	0.3175	0.4567	0.3673

Table 1: Quantitative comparison of linear dimension reduction methods. The best scores for each metric are bolded.

PCA. For PCA, I retained the first two principal components and projected the samples into the 2D subspace they span, which preserves the primary directions of variance. As for hyperparameters, PCA has very few tunable parameters—the main one being the number of principal components. For ease of visualization, I set the number of components to 2

(PC=2) in this experiment. Figure 1a shows that the data are well dispersed, and digits such as 0, 3, 4, 6, and 9 form clearly separated clusters. Quantitatively, PCA achieved an ARI of 0.3614, an NMI of 0.5190, and a Silhouette score of 0.3993. Except for the Silhouette score, PCA obtained the highest values among the three methods. This indicates that PCA effectively separates the digits and maintains a high level of consistency with the true labels. Although its Silhouette value (approximately 0.4) is not the highest, it still suggests reasonably compact and well-separated clusters. This minor difference can be attributed to slight overlaps between neighboring clusters in the 2D embedding, even though the overall structure aligns well with the ground-truth classes.

In addition to the 2D embedding, I plotted the PCA scree plot and a bar chart of the top ten principal components' explained variance, as shown in Figure 2. Both plots reveal that the first three components contribute substantially more variance than the rest, with the third component explaining slightly less variance than the first two but significantly more than the fourth. This supports the observation that the 2D projection loses some discriminative information, which explains why the best ARI achieved by PCA remains moderate (0.3614).

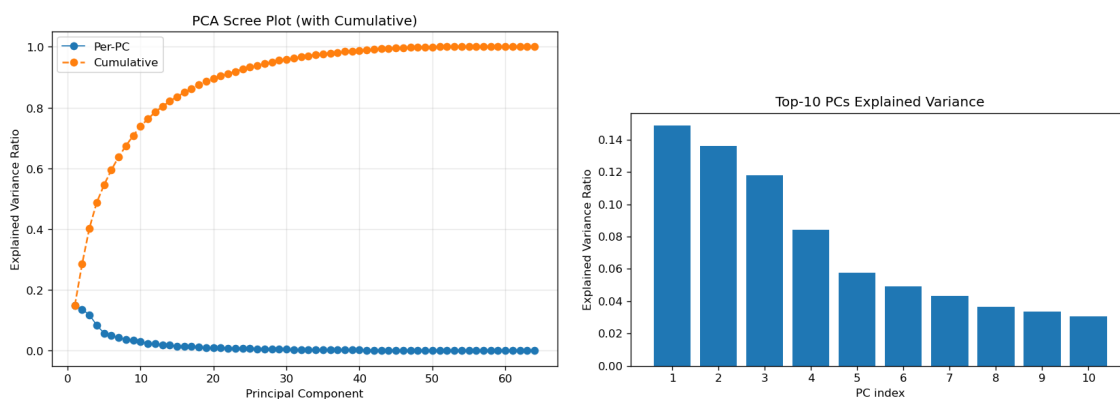


Figure 2: (Left) Scree plot showing the variance explained by each principal component; (Right) bar chart of the top-10 PCs' explained variance ratios.

Furthermore, I visualized the top ten PCA component images (Figure 3), which illustrate the principal modes of variation across digits.

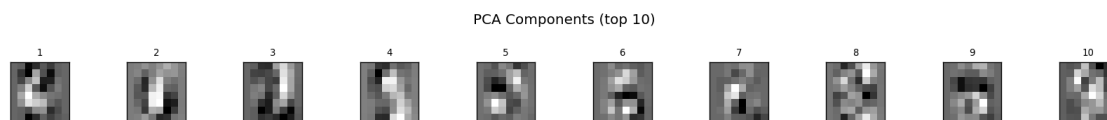


Figure 3: Top ten PCA components visualized as 8×8 basis images.

NMF→PCA. For the NMF method, I first determined the optimal number of components by minimizing the reconstruction error, which yielded $k = 20$. This means the data were decomposed into 20 non-negative basis vectors. The resulting coefficient matrix W was

then compressed into a 2D embedding space using PCA for visualization and comparison purposes, and the final result is shown in Figure 1(b).

As observed in the plot, NMF fails to clearly separate the digits in the 2D space. This is consistent with the quantitative results, where NMF achieved the lowest ARI (0.1588) and NMI (0.2961) among all methods. These findings indicate that the NMF representation captures local, part-based features rather than global discriminative structures, and that the subsequent PCA compression may distort these part-based patterns, reducing the overall interpretability.

On the other hand, NMF obtained the highest Silhouette score (0.4080), slightly higher than PCA. The embedding shows that nearly all samples are concentrated in the positive subspace, with only a small portion extending below zero (no less than -0.2). Even after PCA compression, this non-negativity-induced structure is largely preserved, resulting in an asymmetric distribution that nevertheless exhibits slightly better cluster compactness than PCA. This suggests that some of the features extracted by NMF may be better suited to local clustering in this dataset. In particular, compared to PCA, NMF tends to focus on localized regions of variation rather than global variance directions, which likely contributes to the formation of tighter, more compact clusters.

Similar to the PCA case where excluding the third principal component led to a loss of explanatory power, it is plausible that the 2D projection of NMF also omits important structural information. Specifically, the visualization reveals that the first principal component successfully separates digit “4” along the left margin, while the second principal component distinguishes digits “9” (light blue) and “0” from the rest of the samples, forming a clear boundary near the bottom region of the plot. This indicates that the first PC primarily captures the unique pattern of the digit “4”, whereas the second PC isolates the shapes shared by “0” and “9”. A potential 3D embedding including an additional principal component might further clarify the remaining overlapping clusters visible in the upper-right portion of the 2D space.

Finally, I visualized the NMF basis components, as shown in Figure 4. These components represent localized stroke-like patterns, providing an interpretable decomposition of the digits into additive parts.

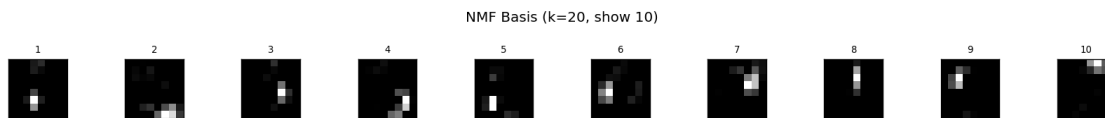


Figure 4: Top ten NMF basis components visualized as 8×8 basis images.

ICA. For the ICA method, I applied FastICA with two components after standardizing the data to ensure consistent feature scaling. As shown in Figure 1(c) and the summary table, ICA exhibits the most balanced overall performance among the three linear methods. Some clusters—notably digits 4, 6, and 0—are well separated and relatively compact, outperforming both PCA and NMF in local structure preservation. However, the remaining digits appear more mixed than in PCA, though less dispersed than in NMF. Visually,

the embedding forms two major groups: one consisting primarily of digits 4, 6, and 0, and another containing the other seven digits.

Quantitatively, ICA’s ARI (0.3175) and NMI (0.4567) values lie between those of PCA and NMF, while its Silhouette score (0.3673) is the lowest among the three. This aligns with the visual interpretation: ICA achieves moderate global separability but weaker overall cluster compactness.

Regarding hyperparameters, ICA provides limited tuning options. Here I set `components=2`, reducing the data to a 2D space, and standardized all features prior to decomposition. Standardization is a conventional preprocessing step for ICA, as its underlying assumption relies on statistical independence between components. Consequently, the resulting embedding appears nearly spherical and centered around the origin—a distribution consistent with ICA’s model assumptions but inconsistent with the inherent non-spherical structure of handwritten digits. This mismatch explains the weaker performance of ICA in this task.

Finally, the ICA component images (Figure 5) reveal contrast-like and edge-detecting patterns that emphasize stroke boundaries, differing from the smoother, global variations captured by PCA and the localized parts extracted by NMF.

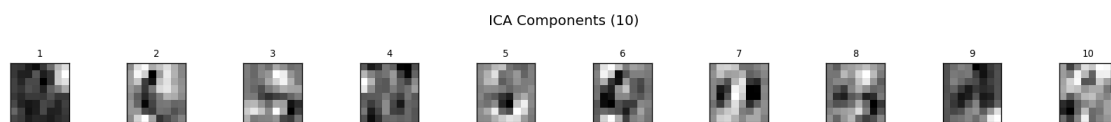


Figure 5: Top ten PCA components visualized as 8×8 basis images.

Overall Discussion. In summary, among the three linear dimension reduction techniques, PCA stands out as the most effective approach. It achieved the highest ARI (0.3614) and NMI (0.5190), and a Silhouette score (0.3993) close to the best. Visually, PCA produced the clearest and most interpretable clusters in the 2D embedding, separating several digits (such as 0, 3, 4, 6, and 9) distinctly. Given that PCA captures global variance directions, its performance would likely improve further in higher-dimensional embeddings, where additional components could better preserve discriminative variance.

1.2 Apply manifold learning approaches.

Methodology Each model was configured with typical hyperparameters and applied to the same standardized input features. The detailed settings are summarized below:

- **Kernel PCA (RBF)** — RBF kernel with $\gamma = 0.102$, chosen automatically based on the median pairwise distance.
- **Spectral Embedding** — Laplacian eigenmap-based manifold embedding with $n_{\text{neighbors}} = 10$.
- **Classical MDS** — Closed-form double-centered Euclidean distance eigendecomposition.

- **Metric MDS** — Stress minimization–based metric scaling via gradient descent.
- **t-SNE** — Perplexity of 30, PCA initialization, adaptive learning rate.
- **UMAP** — $n_{\text{neighbors}} = 15$, $\text{min_dist} = 0.1$.
- **Autoencoder (PyTorch)** — Fully connected symmetric architecture with input 64, hidden layer 32, bottleneck dimension 2, ReLU activations, sigmoid reconstruction, trained for 50 epochs with batch size 128 and learning rate 10^{-3} .

All embeddings were evaluated in 2D using k -means clustering ($k = 10$) with three quantitative metrics: Adjusted Rand Index (ARI), Normalized Mutual Information (NMI), and Silhouette coefficient.

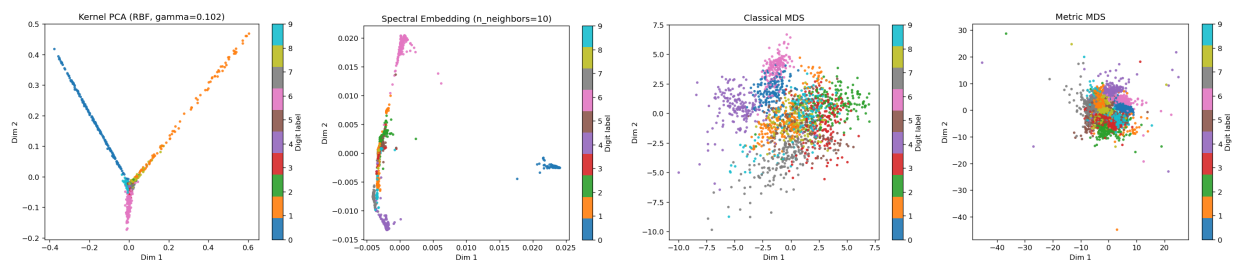


Figure 6: 2D embeddings of the digits dataset using kernel PCA, spectral embedding, classical mds and metric mds.

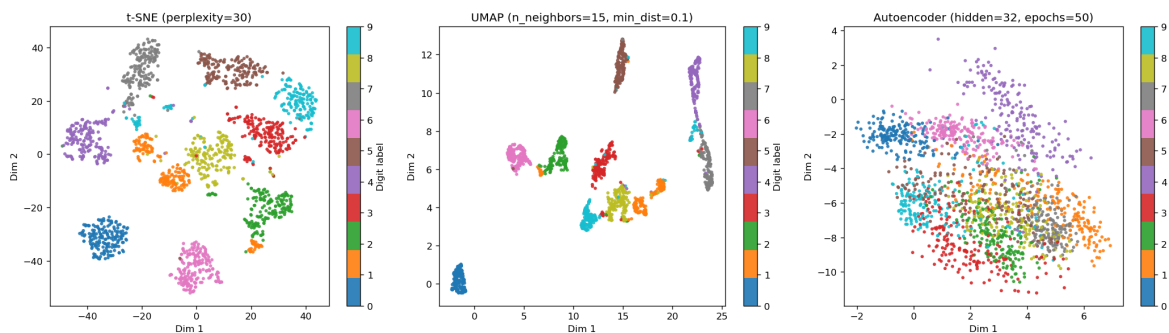


Figure 7: 2D embeddings of the digits dataset using t-SNE, UMAP and Pytorch Autoencoder.

Discussion The quantitative results and visualizations reveal clear distinctions among methods. UMAP and t-SNE perform best overall, yielding the strongest alignment with true digit labels—UMAP achieves the highest ARI (0.866) and NMI (0.891), while t-SNE produces similarly well-separated and visually coherent clusters. Spectral Embedding captures both global and local structures but exhibits partial overlap between classes. Kernel PCA attains the highest Silhouette score (0.754), reflecting locally compact but label-misaligned clusters, whereas Classical and Metric MDS show limited nonlinear capacity. The autoencoder provides a moderate nonlinear representation (ARI 0.35, NMI 0.50), illustrating

Method	ARI	NMI	Silhouette
Kernel PCA	0.0511	0.3357	0.7541
Spectral Embedding	0.4730	0.6681	0.6119
Classical MDS	0.3262	0.4647	0.3771
Metric MDS	0.3461	0.4722	0.3503
t-SNE	0.7695	0.8339	0.5731
UMAP	0.8662	0.8912	0.6866
Autoencoder	0.3476	0.4992	0.3958

Table 2: Quantitative comparison of seven embedding methods on the digits dataset ($k = 10$).

its potential but also the constraint of a shallow 2D bottleneck. Overall, UMAP and t-SNE best uncover the intrinsic manifold structure of handwritten digits, highlighting the strength of nonlinear embedding methods in revealing meaningful low-dimensional organization in high-dimensional data.

1.3 Discussion

Nonlinear manifold methods such as **UMAP** and **t-SNE** outperform linear approaches because they better preserve both local neighborhood relationships and global manifold geometry, thus capturing the intrinsic nonlinear structure of handwritten digits. Among all methods, **UMAP** performs best, achieving the highest ARI (0.866) and NMI (0.891) while forming compact and well-separated clusters that closely align with true digit classes. The most representative visualization is the **UMAP 2D embedding**, which clearly displays ten distinct clusters corresponding to the ten digits.

2 Open-Ended Data Analysis - Breast Cancer gene expression data.

2.1 Preprocessing

After loading the data (shape 445×359), six clinical columns were identified: **subtype**, **er_status**, **pr_status**, **her2_status**, **node**, and **metastasis**, leaving 353 gene expression features. No missing values or zero-variance genes were found, indicating high data quality. Each gene feature was standardized via Z-score normalization. Descriptive statistics of the clinical variables are summarized below:

- **Subtype:** Luminal A (200), Luminal B (106), Basal-like (79), HER2-enriched (53), Normal-like (7).
- **ER-Status:** Positive (339), Negative (100), Performed but Not Available (2), Indeterminate (2), Not Performed (2).

- **PR-Status:** Positive (291), Negative (147), Indeterminate (3), Performed but Not Available (2), Not Performed (2).
- **HER2-Status:** Negative (371), Positive (65), Equivocal (5), Not Available (4).
- **Node:** mean = 0.73, std = 0.87, range = [0, 3].
- **Metastasis:** mean = 0.025, std = 0.155 (mostly non-metastatic samples).

2.2 Methodology

Five dimension reduction methods were applied to the standardized gene expression matrix to obtain two-dimensional embeddings:

1. **Principal Component Analysis (PCA)** — linear orthogonal projection capturing maximal variance.
2. **Non-negative Matrix Factorization (NMF)** — parts-based representation using non-negative constraints (run on non-standardized data to ensure non-negativity).
3. **Spectral Embedding** — manifold learning based on graph Laplacian eigenvectors.
4. **t-SNE** — nonlinear embedding preserving local neighborhood structures.
5. **UMAP** — manifold approximation balancing local and global relationships.

Each embedding was visualized by coloring the points according to all six available clinical variables. However, for brevity and visual clarity, only the results colored by **Subtype** are shown here as representative examples. This selection illustrates overall trends while the complete set of figures (for all six variables across all five methods) is provided in the supplementary materials.

To quantitatively evaluate clustering quality with respect to molecular subtypes, k -means clustering ($k = 5$) was applied to each embedding, and three metrics were computed: Adjusted Rand Index (ARI), Normalized Mutual Information (NMI), and Silhouette coefficient.

2.3 Results

Overall, nonlinear manifold-based methods (t-SNE and UMAP) outperformed linear approaches in terms of ARI and NMI, suggesting that the underlying structure of gene expression data is highly nonlinear. UMAP achieved the highest Silhouette score (0.44), indicating that it produces well-separated clusters with clear boundaries, while t-SNE yielded the best ARI (0.21) and NMI (0.27), showing the strongest alignment with the known PAM50 subtypes. PCA achieved moderate performance, confirming its ability to preserve global variance but limited capacity to capture complex nonlinear relations. NMF performed comparably, producing slightly denser local clusters (reflected in its higher Silhouette score) but weaker subtype separation. Spectral embedding produced the most compact clusters but with limited biological interpretability due to its sensitivity to graph construction.

Method	ARI	NMI	Silhouette
PCA	0.1857	0.2399	0.3311
NMF	0.1712	0.2536	0.3842
Spectral	0.1421	0.2300	0.4020
t-SNE	0.2074	0.2704	0.3984
UMAP	0.2028	0.2780	0.4376

Table 3: Quantitative comparison of five dimension reduction methods on the BRCA dataset. Evaluation is based on clustering alignment with PAM50 subtypes ($k = 5$).

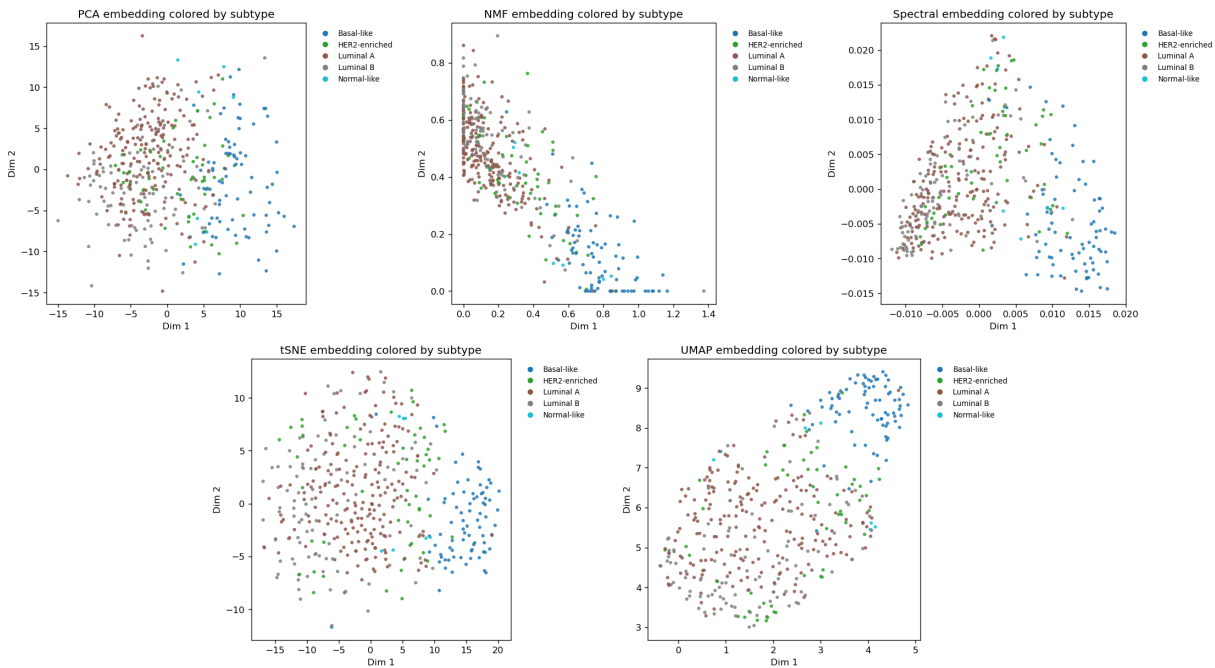


Figure 8: Two-dimensional embeddings of the BRCA gene expression data using five methods (PCA, NMF, Spectral, t-SNE, and UMAP), colored by molecular Subtype. All embeddings were also generated for the remaining five clinical variables, but only Subtype-colored results are displayed here for clarity.

2.4 Discussion

Given that no missing or low-variance genes were present, preprocessing had minimal impact on the raw data distribution. The observed performance differences mainly stem from the intrinsic characteristics of each method. The superior performance of UMAP and t-SNE suggests that manifold learning effectively captures nonlinear relationships among gene expressions that correspond to known molecular subtypes. Future work could extend this analysis by increasing latent dimensionality or incorporating autoencoders for deeper non-linear representations.

A Appendix: Code Implementation

A.1 Problem 1a.

```

1  '''
2  Author: Chuyang Su cs4570@columbia.edu
3  Date: 2025-10-08 17:16:54
4  LastEditTime: 2025-10-08 18:21:50
5  FilePath: /Unsupervised-Learning-Homework/Homework 1/Code/Problem_1_a.py
6  Description:
7      This is the code part of GR5244 Unsupervised Learning Homework 1 Part 1a.
8  '''
9  import argparse
10 import os
11 from pathlib import Path
12
13 import numpy as np
14 import matplotlib.pyplot as plt
15
16 from sklearn.datasets import load_digits
17 from sklearn.decomposition import PCA, NMF, FastICA
18 from sklearn.cluster import KMeans
19 from sklearn.metrics import adjusted_rand_score, normalized_mutual_info_score,
20     silhouette_score
21 from sklearn.preprocessing import StandardScaler
22 from sklearn.pipeline import make_pipeline
23 import csv
24
25 # -----
26 # Utils
27 # -----
28 def ensure_dir(p):
29     Path(p).mkdir(parents=True, exist_ok=True)
30
31 def plot_embedding(X2d, y, title, outpath):
32     plt.figure(figsize=(6, 5), dpi=120)
33     scatter = plt.scatter(X2d[:, 0], X2d[:, 1], c=y, s=12, cmap='tab10', alpha
34         =0.9, edgecolors='none')
35
36     plt.title(title)
37     plt.xlabel("Dim 1")
38     plt.ylabel("Dim 2")
39     cbar = plt.colorbar(scatter, ticks=range(10))
40     cbar.set_label("Digit label")
41     plt.tight_layout()
42     plt.savefig(outpath, bbox_inches='tight')
43     plt.close()
44
45 def plot_components(components, img_shape, n_show, title, outpath):
46     n = min(n_show, components.shape[0])
47     cols = 10 if n >= 10 else n
48     rows = int(np.ceil(n / cols))

```

```

48 plt.figure(figsize=(1.4*cols, 1.4*rows), dpi=120)
49 for i in range(n):
50     ax = plt.subplot(rows, cols, i + 1)
51     ax.imshow(components[i].reshape(img_shape), cmap='gray')
52     ax.set_xticks([])
53     ax.set_yticks([])
54     ax.set_title(f"{i+1}", fontsize=8)
55 plt.suptitle(title, y=1.02)
56 plt.tight_layout()
57 plt.savefig(outpath, bbox_inches='tight')
58 plt.close()
59
60 def plot_pca_scree(explained_var_ratio, outpath):
61     """Scree plot: per-PC variance ratio + cumulative curve"""
62     import numpy as np
63     import matplotlib.pyplot as plt
64
65     r = np.array(explained_var_ratio)
66     cum = np.cumsum(r)
67
68     plt.figure(figsize=(7, 5), dpi=120)
69     #
70     plt.plot(np.arange(1, len(r)+1), r, marker='o', linewidth=1)
71     #
72     plt.plot(np.arange(1, len(r)+1), cum, marker='o', linestyle='--')
73     plt.xlabel("Principal Component")
74     plt.ylabel("Explained Variance Ratio")
75     plt.title("PCA Scree Plot (with Cumulative)")
76     plt.legend(["Per-PC", "Cumulative"])
77     plt.grid(alpha=0.3)
78     plt.tight_layout()
79     plt.savefig(str(outpath), bbox_inches='tight')
80     plt.close()
81
82
83 def plot_pca_topk_bar(explained_var_ratio, k, outpath):
84     """Bar chart for top-k PCs' explained variance ratio"""
85     import numpy as np
86     import matplotlib.pyplot as plt
87
88     r = np.array(explained_var_ratio)[:k]
89     idx = np.arange(1, len(r)+1)
90
91     plt.figure(figsize=(7, 4), dpi=120)
92     plt.bar(idx, r)
93     plt.xlabel("PC index")
94     plt.ylabel("Explained Variance Ratio")
95     plt.title(f"Top-{k} PCs Explained Variance")
96     plt.xticks(idx)
97     plt.tight_layout()
98     plt.savefig(str(outpath), bbox_inches='tight')
99     plt.close()
100
101

```

```

102
103 def kmeans_scores(X2d, y, seed):
104     km = KMeans(n_clusters=10, n_init='auto', random_state=seed)
105     labels = km.fit_predict(X2d)
106     ari = adjusted_rand_score(y, labels)
107     nmi = normalized_mutual_info_score(y, labels)
108     sil = silhouette_score(X2d, labels)
109     return ari, nmi, sil
110
111
112 # -----
113 # Main pipeline for 1(a)
114 # -----
115 def run(seed=0, nmf_ks=(10, 15, 20), outdir="Homework 1/Latex"):
116     ensure_dir(outdir)
117
118     # Load data
119     digits = load_digits()
120     X = digits.data.astype(float)
121     y = digits.target
122     img_shape = (8, 8)
123
124     # -----
125     # PCA (2D embedding + components)
126     # -----
127     pca_2 = PCA(n_components=2, random_state=seed)
128     X_pca_2 = pca_2.fit_transform(X)
129     plot_embedding(X_pca_2, y, "PCA (2D)", f"{outdir}/pca_2d.png")
130
131     # For components visualization, use more PCs (e.g., 10)
132     pca_10 = PCA(n_components=10, random_state=seed).fit(X)
133     plot_components(pca_10.components_, img_shape, n_show=10,
134                    title="PCA Components (top 10)", outpath=f"{outdir}/
135                                                                pca_components.png"
136                    )
137
138     pca_ari, pca_nmi, pca_sil = kmeans_scores(X_pca_2, y, seed)
139
140     pca_full = PCA(n_components=min(X.shape), random_state=seed).fit(X)
141     explained = pca_full.explained_variance_ratio_
142
143     plot_pca_scree(explained, Path(outdir) / "pca_scree.png")
144     plot_pca_topk_bar(explained, k=10, outpath=Path(outdir) / "pca_top10_var.
145                                                                png")
146
147     # -----
148     # NMF (grid over k), then 2D via PCA-on-W for visualization
149     # -----
150     X_nonneg = X - X.min() if X.min() < 0 else X
151     best_nmf = None
152     best_rec = np.inf
153     best_k = None

```

```

153     for k in nmf_ks:
154         nmf = NMF(n_components=k, init='nndsvda', random_state=seed, max_iter=
                                1000)
155         W = nmf.fit_transform(X_nonneg)
156         rec = nmf.reconstruction_err_
157         if rec < best_rec:
158             best_rec = rec
159             best_nmf = nmf
160             best_k = k
161
162     # Use the best NMF
163     W_best = best_nmf.transform(X_nonneg) # (n_samples, best_k)
164     # Reduce W to 2D for visualization
165     nmf_to2 = PCA(n_components=2, random_state=seed)
166     X_nmf_2 = nmf_to2.fit_transform(W_best)
167     plot_embedding(X_nmf_2, y, f"NMF -> PCA (2D) [k={best_k}]", f"{outdir}/
                                nmf_2d.png")
168
169     # Visualize NMF basis (H)
170     plot_components(best_nmf.components_, img_shape, n_show=10,
171                    title=f"NMF Basis (k={best_k}, show 10)", outpath=f"{
                                outdir}/
                                nmf_components.png"
                                )
172
173     nmf_ari, nmf_nmi, nmf_sil = kmeans_scores(X_nmf_2, y, seed)
174
175     # -----
176     # ICA (2D embedding + components)
177     # -----
178     ica_pipeline = make_pipeline(StandardScaler(with_std=True), FastICA(
                                n_components=2, random_state=seed,
                                max_iter=1000))
179     X_ica_2 = ica_pipeline.fit_transform(X)
180     plot_embedding(X_ica_2, y, "ICA (2D)", f"{outdir}/ica_2d.png")
181
182     # For components display, fit a separate ICA with more comps (e.g., 10) on
        standardized X
183     scaler = StandardScaler(with_std=True)
184     X_std = scaler.fit_transform(X)
185     ica_10 = FastICA(n_components=10, random_state=seed, max_iter=1000).fit(
                                X_std)
186     plot_components(ica_10.mixing_.T, img_shape, n_show=10, # mixing_.T
                                component ""images
187                    title="ICA Components (10)", outpath=f"{outdir}/
                                ica_components.png"
                                )
188
189     ica_ari, ica_nmi, ica_sil = kmeans_scores(X_ica_2, y, seed)
190
191     # -----
192     # Save metrics
193     # -----
194     metrics_path = f"{outdir}/part1a_metrics.csv"

```

```

195 with open(metrics_path, "w", newline="") as f:
196     writer = csv.writer(f)
197     writer.writerow(["Method", "Params", "ARI", "NMI", "Silhouette", "
                        Notes"])
198     writer.writerow(["PCA", "n_components=2", f"{pca_ari:.4f}", f"{pca_nmi
                        :.4f}", f"{pca_sil:.4f}", "2D
                        embedding"])
199     writer.writerow(["NMF -> PCA", f"k={best_k}, then 2D PCA", f"{nmf_ari
                        :.4f}", f"{nmf_nmi:.4f}", f"{
                        nmf_sil:.4f}",
200                     f"best reconstruction k among {list(nmf_ks)} (err={
                        best_rec:.4f})
                        "] )
201     writer.writerow(["ICA", "n_components=2 (with standardization)", f"{
                        ica_ari:.4f}", f"{ica_nmi:.4f}"
                        , f"{ica_sil:.4f}", "2D
                        embedding"])
202
203 # Also print a neat summary
204 print("\n=== Part 1(a) - Linear Methods on Digits ===")
205 print(f"[PCA]          ARI={pca_ari:.4f}  NMI={pca_nmi:.4f}  Silhouette={
                        pca_sil:.4f}")
206 print(f"[NMF -> PCA]   ARI={nmf_ari:.4f}  NMI={nmf_nmi:.4f}  Silhouette={
                        nmf_sil:.4f}  (best k={best_k},
                        recon_err={best_rec:.4f})")
207 print(f"[ICA]          ARI={ica_ari:.4f}  NMI={ica_nmi:.4f}  Silhouette={
                        ica_sil:.4f}")
208 print(f"\nSaved figures and metrics to: {outdir}/")
209 print("Figures:")
210 print(" - pca_2d.png, pca_components.png")
211 print(" - nmf_2d.png, nmf_components.png")
212 print(" - ica_2d.png, ica_components.png")
213 print("Table:")
214 print(f" - {Path(metrics_path).name}")
215
216
217 # -----
218 # Entry
219 # -----
220 if __name__ == "__main__":
221     parser = argparse.ArgumentParser()
222     parser.add_argument("--seed", type=int, default=25)
223     parser.add_argument("--outdir", type=str, default="Homework 1/Latex/
                        Results/Problem_1_a")
224     parser.add_argument("--nmf_ks", type=int, nargs="+", default=[10, 15, 20])
225     args = parser.parse_args()
226     run(seed=args.seed, nmf_ks=tuple(args.nmf_ks), outdir=args.outdir)

```

A.2 Problem 1b.

```

1 #!/usr/bin/env python3
2 # -*- coding: utf-8 -*-

```

```

3 """
4 Author: Chuyang Su <cs4570@columbia.edu>
5 Date: 2025-10-10
6 FilePath: /Unsupervised-Learning-Homework/Homework 1/Code/Problem_1_b.py
7 Description:
8     GR5244 HW1 Part 1(b): Manifold learning on sklearn digits (n=1797, p=64).
9     Methods: Kernel PCA, Spectral Embedding, Classical MDS, Metric MDS, t-
10              SNE, UMAP, Autoencoder.
11 """
12 from pathlib import Path
13 import argparse
14 import numpy as np
15 import matplotlib.pyplot as plt
16
17 from sklearn.datasets import load_digits
18 from sklearn.preprocessing import StandardScaler, MinMaxScaler
19 from sklearn.decomposition import KernelPCA
20 from sklearn.manifold import SpectralEmbedding, MDS, TSNE
21 from sklearn.metrics import adjusted_rand_score, normalized_mutual_info_score,
22     silhouette_score, pairwise_distances
23 from sklearn.cluster import KMeans
24 import umap
25
26 import torch
27 import torch.nn as nn
28 import torch.optim as optim
29 from torch.utils.data import TensorDataset, DataLoader
30 # -----
31 # Utils
32 # -----
33 def ensure_dir(p):
34     Path(p).mkdir(parents=True, exist_ok=True)
35
36 def set_seed(seed: int):
37     import random, os
38     random.seed(seed); np.random.seed(seed)
39     torch.manual_seed(seed); torch.cuda.manual_seed_all(seed)
40     torch.backends.cudnn.deterministic = True
41     torch.backends.cudnn.benchmark = False
42
43 def get_device():
44     return torch.device("cuda" if torch.cuda.is_available() else "cpu")
45
46 def plot_embedding(X2d, y, title, outpath):
47     """Simple 2D scatter with digits colormap; no axes ticks; safe colormap
48         API."""
49     cmap = plt.colormaps.get_cmap('tab10')
50     plt.figure(figsize=(6, 5), dpi=120)
51     sc = plt.scatter(X2d[:, 0], X2d[:, 1], c=y, s=12, cmap=cmap, alpha=0.9,
52         edgecolors='none')
53     plt.title(title)
54     plt.xlabel("Dim 1"); plt.ylabel("Dim 2")

```

```

53     cbar = plt.colorbar(sc, ticks=range(10))
54     cbar.set_label("Digit label")
55     plt.tight_layout()
56     plt.savefig(outpath, bbox_inches='tight')
57     plt.close()
58
59
60 def kmeans_scores(X2d, y, seed, k=10):
61     km = KMeans(n_clusters=k, random_state=seed, n_init=10)
62     labels = km.fit_predict(X2d)
63     ari = adjusted_rand_score(y, labels)
64     nmi = normalized_mutual_info_score(y, labels)
65     sil = silhouette_score(X2d, labels)
66     return ari, nmi, sil
67
68
69 def classical_mds(X, n_components=2):
70     """
71     Classical MDS (Torgerson/Gower): eigendecomposition of double-centered
72     squared-distance matrix. Returns low-d coords (can be rotated/reflected).
73     """
74     D = pairwise_distances(X, metric='euclidean')
75     D2 = D ** 2
76     n = D2.shape[0]
77     J = np.eye(n) - np.ones((n, n)) / n
78     B = -0.5 * J @ D2 @ J
79     # Eigen-decomposition
80     evals, evects = np.linalg.eigh(B)
81     # Take top components
82     idx = np.argsort(evals)[::-1]
83     evals = evals[idx]; evects = evects[:, idx]
84     # Keep only positive eigenvalues
85     pos = evals > 0
86     evals = evals[pos]; evects = evects[:, pos]
87     evals_k = evals[:n_components]
88     evects_k = evects[:, :n_components]
89     X_emb = evects_k * np.sqrt(np.maximum(evals_k, 0))
90     return X_emb
91
92
93 def auto_gamma_rbf(X):
94     """
95     Heuristic gamma for RBF kernel (Kernel PCA): gamma = 1 / median(pairwise
96     distance).
97     Scale-robust; avoids manual guesswork when not provided.
98     """
99     d = pairwise_distances(X, metric='euclidean')
100     med = np.median(d)
101     if med <= 0:
102         return 1.0
103     return 1.0 / med
104
105 # Autoencoder

```

```

106 class AE(nn.Module):
107     def __init__(self, input_dim: int, hidden: int = 32, bottleneck: int = 2):
108         super().__init__()
109         self.encoder = nn.Sequential(
110             nn.Linear(input_dim, hidden),
111             nn.ReLU(),
112             nn.Linear(hidden, bottleneck) #
113         )
114         self.decoder = nn.Sequential(
115             nn.Linear(bottleneck, hidden),
116             nn.ReLU(),
117             nn.Linear(hidden, input_dim),
118             nn.Sigmoid() # [0,1]
119         )
120
121     def forward(self, x):
122         z = self.encoder(x)
123         xhat = self.decoder(z)
124         return xhat, z
125
126
127 @torch.no_grad()
128 def encode_dataset(model: AE, loader: DataLoader, device):
129     model.eval()
130     zs = []
131     for (xb,) in loader:
132         xb = xb.to(device)
133         _, z = model(xb)
134         zs.append(z.cpu().numpy())
135     return np.concatenate(zs, axis=0)
136
137
138 def train_autoencoder_pytorch(
139     X_minmax: np.ndarray,
140     seed: int = 0,
141     hidden: int = 32,
142     bottleneck: int = 2,
143     epochs: int = 50,
144     batch_size: int = 128,
145     lr: float = 1e-3,
146 ):
147     set_seed(seed)
148     device = get_device()
149
150     X_tensor = torch.tensor(X_minmax, dtype=torch.float32)
151     ds = TensorDataset(X_tensor)
152     dl = DataLoader(ds, batch_size=batch_size, shuffle=True, drop_last=False)
153
154     model = AE(input_dim=X_minmax.shape[1], hidden=hidden, bottleneck=
                    bottleneck).to(device)
155     opt = optim.Adam(model.parameters(), lr=lr)
156     crit = nn.MSELoss()
157     model.train()
158     for _ in range(epochs):

```



```

159         for (xb,) in dl:
160             xb = xb.to(device)
161             xhat, _ = model(xb)
162             loss = crit(xhat, xb)
163             opt.zero_grad()
164             loss.backward()
165             opt.step()
166
167         #      2D
168         #      DataLoader
169         Z = encode_dataset(model, DataLoader(ds, batch_size=1024, shuffle=False),
170                                device)
171
172         return Z
173
174 # -----
175 # Main pipeline for 1(b)
176 # -----
177 def run(seed=25,
178         outdir="Homework 1/Latex/Results/Problem_1_b",
179         kpca_gamma=None,
180         spectral_n_neighbors=10,
181         tsne_perplexity=30,
182         umap_n_neighbors=15,
183         umap_min_dist=0.1,
184         ae_hidden=32,
185         ae_epochs=50,
186         ae_batch=128,
187         ae_lr=1e-3):
188     set_seed(seed)
189     ensure_dir(outdir)
190
191     # Load data
192     digits = load_digits()
193     X = digits.data.astype(float)          # (1797, 64)
194     y = digits.target                      # (1797,)
195
196     # Standardize (common for manifold learning)
197     X_std = StandardScaler().fit_transform(X)
198
199     # 1) Kernel PCA (RBF)
200     gamma = kpca_gamma if kpca_gamma is not None else auto_gamma_rbf(X_std)
201     kpca = KernelPCA(n_components=2, kernel='rbf', gamma=gamma,
202                      fit_inverse_transform=False,
203                      random_state=seed)
204
205     X_kpca = kpca.fit_transform(X_std)
206     plot_embedding(X_kpca, y, f"Kernel PCA (RBF, gamma={gamma:.3g})", f"{
207                                outdir}/kpca_2d.png")
208     kpca_ari, kpca_nmi, kpca_sil = kmeans_scores(X_kpca, y, seed)
209
210     # 2) Spectral Embedding
211     se = SpectralEmbedding(n_components=2, n_neighbors=spectral_n_neighbors,
212                            random_state=seed)
213
214     X_se = se.fit_transform(X_std)

```

```

208 plot_embedding(X_se, y, f"Spectral Embedding (n_neighbors={
                                spectral_n_neighbors})", f"{outdir}
                                /spectral_2d.png")
209 se_ari, se_nmi, se_sil = kmeans_scores(X_se, y, seed)
210
211 # 3) Classical MDS (closed-form)
212 X_cmds = classical_mds(X_std, n_components=2)
213 plot_embedding(X_cmds, y, "Classical MDS", f"{outdir}/classical_mds_2d.png
                                ")
214 cmds_ari, cmds_nmi, cmds_sil = kmeans_scores(X_cmds, y, seed)
215
216 # 4) Metric MDS (sklearn MDS, metric=True)
217 mds_metric = MDS(n_components=2, metric=True, normalized_stress='auto',
                                random_state=seed)
218 X_mds = mds_metric.fit_transform(X_std)
219 plot_embedding(X_mds, y, "Metric MDS", f"{outdir}/metric_mds_2d.png")
220 mds_ari, mds_nmi, mds_sil = kmeans_scores(X_mds, y, seed)
221
222 # 5) t-SNE
223 tsne = TSNE(n_components=2, perplexity=tsne_perplexity, learning_rate='
                                auto', init='pca', random_state=
                                seed)
224 X_tsne = tsne.fit_transform(X_std)
225 plot_embedding(X_tsne, y, f"t-SNE (perplexity={tsne_perplexity})", f"{
                                outdir}/tsne_2d.png")
226 tsne_ari, tsne_nmi, tsne_sil = kmeans_scores(X_tsne, y, seed)
227
228 # 6) UMAP
229 umap_model = umap.UMAP(n_neighbors=umap_n_neighbors, min_dist=
                                umap_min_dist, n_components=2,
                                random_state=seed)
230 X_umap = umap_model.fit_transform(X_std)
231 plot_embedding(X_umap, y, f"UMAP (n_neighbors={umap_n_neighbors}, min_dist
                                ={umap_min_dist})", f"{outdir}/
                                umap_2d.png")
232 umap_ari, umap_nmi, umap_sil = kmeans_scores(X_umap, y, seed)
233
234 # AutoencoderPyTorch
235 X_minmax = MinMaxScaler().fit_transform(X)
236 X_ae = train_autoencoder_pytorch(
237     X_minmax,
238     seed=seed,
239     hidden=ae_hidden,
240     bottleneck=2,
241     epochs=ae_epochs,
242     batch_size=ae_batch,
243     lr=ae_lr,
244 )
245 plot_embedding(X_ae, y, f"Autoencoder (hidden={ae_hidden}, epochs={
                                ae_epochs})", f"{outdir}/
                                autoencoder_2d.png")
246 ae_ari, ae_nmi, ae_sil = kmeans_scores(X_ae, y, seed)
247
248 # Save metrics

```

```

249 import csv
250 metrics_path = f"{outdir}/part1b_metrics.csv"
251 with open(metrics_path, "w", newline="") as f:
252     writer = csv.writer(f)
253     writer.writerow(["Method", "Params", "ARI", "NMI", "Silhouette", "
                        Notes"])
254     writer.writerow(["Kernel PCA", f"rbf, gamma={gamma:.3g}", f"{kpca_ari
                        :.4f}", f"{kpca_nmi:.4f}", f"{
                        kpca_sil:.4f}", "2D embedding"
                        ])
255     writer.writerow(["Spectral", f"n_neighbors={spectral_n_neighbors}", f"
                        {se_ari:.4f}", f"{se_nmi:.4f}",
                        f"{se_sil:.4f}", "2D embedding
                        "])
256     writer.writerow(["Classical MDS", "-", f"{cmds_ari:.4f}", f"{cmds_nmi
                        :.4f}", f"{cmds_sil:.4f}", "
                        closed-form (eigendecomposition
                        )"])
257     writer.writerow(["Metric MDS", "sklearn MDS(metric=True)", f"{mds_ari
                        :.4f}", f"{mds_nmi:.4f}", f"{
                        mds_sil:.4f}", "stress
                        minimization"])
258     writer.writerow(["t-SNE", f"perplexity={tsne_perplexity}", f"{tsne_ari
                        :.4f}", f"{tsne_nmi:.4f}", f"{
                        tsne_sil:.4f}", "2D embedding"
                        ])
259     writer.writerow(["UMAP", f"n_neighbors={umap_n_neighbors}, min_dist={
                        umap_min_dist}", f"{umap_ari:.
                        4f}", f"{umap_nmi:.4f}", f"{
                        umap_sil:.4f}", "2D embedding"
                        ])
260     writer.writerow(["Autoencoder", f"hidden={ae_hidden}, epochs={ae_epochs
                        }, batch={ae_batch}, lr={ae_lr:
                        g}", f"{ae_ari:.4f}", f"{ae_nmi
                        :.4f}", f"{ae_sil:.4f}", "2D
                        bottleneck (PyTorch)"])
261
262 # Print a neat summary
263 print("\n=== Part 1(b) - Manifold Methods on Digits ===")
264 print(f"[Kernel PCA]      ARI={kpca_ari:.4f}  NMI={kpca_nmi:.4f}  Silhouette
        ={{kpca_sil:.4f}}  (gamma={{gamma:.3g}}
        )")
265 print(f"[Spectral]      ARI={{se_ari:.4f}}  NMI={{se_nmi:.4f}}  Silhouette={{
        se_sil:.4f}}  (n_neighbors={{
        spectral_n_neighbors}})")
266 print(f"[Classical MDS] ARI={{cmds_ari:.4f}}  NMI={{cmds_nmi:.4f}}  Silhouette
        ={{cmds_sil:.4f}}")
267 print(f"[Metric MDS]   ARI={{mds_ari:.4f}}  NMI={{mds_nmi:.4f}}  Silhouette={{
        mds_sil:.4f}}")
268 print(f"[t-SNE]        ARI={{tsne_ari:.4f}}  NMI={{tsne_nmi:.4f}}  Silhouette
        ={{tsne_sil:.4f}}  (perplexity={{
        tsne_perplexity}})")
269 print(f"[UMAP]          ARI={{umap_ari:.4f}}  NMI={{umap_nmi:.4f}}  Silhouette
        ={{umap_sil:.4f}}  (n_neighbors={{

```

```

umap_n_neighbors}, min_dist={
umap_min_dist}))")
270 print(f"\nSaved figures and metrics to: {outdir}/")
271 print("Figures:")
272 print(" - kpca_2d.png, spectral_2d.png, classical_mds_2d.png,
metric_mds_2d.png, tsne_2d.png,
umap_2d.png")

273 print("Table:")
274 print(f" - {Path(metrics_path).name}")
275 print(f"[Autoencoder]   ARI={ae_ari:.4f}   NMI={ae_nmi:.4f}   Silhouette={
ae_sil:.4f}   (hidden={ae_hidden},
epochs={ae_epochs})")

276
277 # -----
278 # Entry
279 # -----
280 if __name__ == "__main__":
281     parser = argparse.ArgumentParser()
282     parser.add_argument("--seed", type=int, default=0, help="random seed")
283     parser.add_argument("--outdir", type=str, default="Homework 1/Latex/
Results/Problem_1_b", help="output
directory")
284     parser.add_argument("--kpca_gamma", type=float, default=None, help="RBF
gamma for Kernel PCA (auto if None)
")
285     parser.add_argument("--spectral_n_neighbors", type=int, default=10, help="
n_neighbors for Spectral Embedding"
)
286     parser.add_argument("--tsne_perplexity", type=float, default=30, help="
perplexity for t-SNE")
287     parser.add_argument("--umap_n_neighbors", type=int, default=15, help="
n_neighbors for UMAP")
288     parser.add_argument("--umap_min_dist", type=float, default=0.1, help="
min_dist for UMAP")
289     parser.add_argument("--ae_hidden", type=int, default=32)
290     parser.add_argument("--ae_epochs", type=int, default=50)
291     parser.add_argument("--ae_batch", type=int, default=128)
292     parser.add_argument("--ae_lr", type=float, default=1e-3)
293     args = parser.parse_args()
294
295     run(seed=args.seed,
296         outdir=args.outdir,
297         kpca_gamma=args.kpca_gamma,
298         spectral_n_neighbors=args.spectral_n_neighbors,
299         tsne_perplexity=args.tsne_perplexity,
300         umap_n_neighbors=args.umap_n_neighbors,
301         umap_min_dist=args.umap_min_dist,
302         ae_hidden=args.ae_hidden,
303         ae_epochs=args.ae_epochs,
304         ae_batch=args.ae_batch,
305         ae_lr=args.ae_lr)

```

A.3 Problem 2.

```

1  '''
2  Author: Chuyang Su cs4570@columbia.edu
3  Date: 2025-10-09 13:41:38
4  LastEditors: Please set LastEditors
5  LastEditTime: 2025-10-09 21:33:31
6  FilePath: /Unsupervised-Learning-Homework/Homework 1/Code/Problem_2.py
7  Description:
8      This data set consists of gene expression measurements for n = 445 breast
          cancer tumors and p = 353 genes
          taken from The Cancer Genome Atlas(
          TCGA).
9      This subset of genes was selected based on whether they contain known
          somatic mutations in cancer.
10     Additionally, this data contains clinical data on the
11         (i) Subtype (denotes 5 PAM50 subtypes including Basal-like, Luminal A,
          Luminal B, HER2-enriched, and
          Normal-like),
12         (ii) ER-Status(estrogen-receptor status),
13         (iii) PR-Status (progesterone-receptor status),
14         (iv) HER2-Status (human epidermal growth factor receptor 2 status),
15         (v) Node (number of lymph nodes involved), and (vi)Metastasis (
          indicator for whether the
          cancer has metastasized).
16  '''
17  import pandas as pd
18  import numpy as np
19  from sklearn.preprocessing import StandardScaler
20  from pathlib import Path
21  import os
22  import matplotlib.pyplot as plt
23  from sklearn.decomposition import PCA, NMF
24  from sklearn.manifold import TSNE, SpectralEmbedding
25  from sklearn.cluster import KMeans
26  from sklearn.metrics import adjusted_rand_score, normalized_mutual_info_score,
          silhouette_score
27  import csv
28  import umap
29
30  # === Embedding and evaluation pipeline (UMAP included) ===
31  def run_embedding(X, method='pca', n_components=2, seed=0, **kwargs):
32      """
33      Apply a dimension reduction method and return 2D embedding.
34      Assumes UMAP is available (imported at top).
35      """
36      method = method.lower()
37      if method == 'pca':
38          model = PCA(n_components=n_components, random_state=seed)
39          name = "PCA"
40      elif method == 'nmf':
41          model = NMF(n_components=n_components, init='nndsvda', random_state=
              seed, max_iter=1000)
42          name = "NMF"

```

```

43     elif method == 'tsne':
44         model = TSNE(n_components=n_components, random_state=seed, **kwargs)
45         name = "tSNE"
46     elif method == 'umap':
47         model = umap.UMAP(n_components=n_components, random_state=seed, **
48                             kwargs)
49         name = "UMAP"
50     elif method == 'spectral':
51         model = SpectralEmbedding(n_components=n_components, random_state=seed
52                                   )
53         name = "Spectral"
54     else:
55         raise ValueError(f"Unknown method: {method}")
56     X_embedded = model.fit_transform(X)
57     return X_embedded, name
58
59 def _encode_categories(series):
60     """Encode categorical labels to integer indices and return (indices,
61                                     mapping)."""
62     cats = series.astype(str).fillna("NA").values
63     uniq = sorted(np.unique(cats))
64     mapping = {c: i for i, c in enumerate(uniq)}
65     idx = np.array([mapping[c] for c in cats], dtype=int)
66     return idx, mapping
67
68 def plot_embedding(X2d, meta, hue_col, title, outpath):
69     """Plot and save a 2D embedding colored by a clinical variable."""
70     outpath = Path(outpath)
71     outpath.parent.mkdir(parents=True, exist_ok=True)
72
73     labels = meta[hue_col].astype(str).fillna("NA")
74     idx, mapping = _encode_categories(labels)
75     cmap = plt.cm.get_cmap('tab10', len(mapping))
76
77     plt.figure(figsize=(6.4, 5.2), dpi=120)
78     sc = plt.scatter(X2d[:, 0], X2d[:, 1], c=idx, s=14, cmap=cmap, alpha=0.85,
79                     edgecolors='none')
80     plt.title(f"{title} embedding colored by {hue_col}")
81     plt.xlabel("Dim 1"); plt.ylabel("Dim 2")
82
83     # manual legend
84     handles = [plt.Line2D([0], [0], marker='o', linestyle='', markersize=6,
85                           markerfacecolor=cmap(i), markeredgecolor='none')
86               for i in range(len(mapping))]
87     labels_sorted = list(mapping.keys())
88     plt.legend(handles, labels_sorted, bbox_to_anchor=(1.02, 1), loc='upper
89               left', fontsize=8, frameon=False)
90
91     plt.tight_layout()
92     plt.savefig(outpath, bbox_inches='tight')
93     plt.close()
94
95 def evaluate_embedding(X2d, subtype_series, seed=0):

```

```

92     """KMeans(k=5) on the embedding; compute ARI/NMI/Silhouette using Subtype
93         as reference."""
94     mask = subtype_series.notna()
95     X_eval = X2d[mask.values]
96     y_ref = subtype_series[mask].astype(str).values
97
98     km = KMeans(n_clusters=5, random_state=seed, n_init=10)
99     pred = km.fit_predict(X_eval)
100
101     ari = adjusted_rand_score(y_ref, pred)
102     nmi = normalized_mutual_info_score(y_ref, pred)
103     sil = silhouette_score(X_eval, pred)
104     return ari, nmi, sil
105
106 def run_problem2_pipeline(
107     data_path="Homework 1/Data/BRCA_data.csv",
108     outdir="Homework 1/Latex/Results/Problem_2",
109     seed=0
110 ):
111     outdir = Path(outdir)
112     outdir.mkdir(parents=True, exist_ok=True)
113
114     # 1) Preprocess (writes TXT report under outdir)
115     X, y_clinical, gene_cols = load_and_preprocess_brca(
116         data_path=data_path,
117         report_path=str(outdir / "preprocessing_report.txt")
118     )
119
120     # 2) Methods to run (UMAP included by assumption)
121     methods = [
122         ("pca", dict()),
123         ("nmf", dict()),
124         ("spectral", dict()),
125         ("tsne", dict(perplexity=30, learning_rate='auto', init='pca')),
126         ("umap", dict(n_neighbors=15, min_dist=0.1)),
127     ]
128
129     # 3) Clinical hues to color by (only those that exist)
130     hues = [c for c in ["subtype", "er_status", "pr_status", "her2_status"] if
131              c in y_clinical.columns]
132
133     # 4) Run, plot, evaluate
134     metrics_rows = []
135     best_by_ari = (None, -1.0)
136
137     for m, kwargs in methods:
138         if m == "nmf":
139             # NMF
140             X_nonneg = np.maximum(X, 0) #
141             X2d, name = run_embedding(X_nonneg, method=m, seed=seed, **kwargs)
142         else:
143             X2d, name = run_embedding(X, method=m, seed=seed, **kwargs)

```

```

144     # Evaluation (Subtype as reference, if present)
145     if "subtype" in y_clinical.columns:
146         ari, nmi, sil = evaluate_embedding(X2d, y_clinical["subtype"],
147                                           seed=seed)
148     else:
149         ari = nmi = sil = np.nan
150
151     metrics_rows.append([name, f"{ari:.4f}", f"{nmi:.4f}", f"{sil:.4f}"])
152     if not np.isnan(ari) and ari > best_by_ari[1]:
153         best_by_ari = (name, ari)
154
155     # Save embedding as CSV (optional, useful for debugging/report)
156     emb_path = outdir / f"{name}_embedding.csv"
157     np.savetxt(emb_path, X2d, delimiter=",")
158
159     # Plots colored by clinical variables
160     for hue in hues:
161         fig_path = outdir / f"{name}_{hue}.png"
162         plot_embedding(X2d, y_clinical, hue_col=hue, title=name, outpath=
163                       fig_path)
164
165     # 5) Save metrics summary
166     metrics_path = outdir / "metrics_summary.csv"
167     with open(metrics_path, "w", newline="") as f:
168         writer = csv.writer(f)
169         writer.writerow(["Method", "ARI (Subtype)", "NMI (Subtype)", "
170                         Silhouette (KMeans=5)"])
171         writer.writerows(metrics_rows)
172
173     # Console summary
174     print("\n=== Problem 2: embedding metrics (Subtype as reference) ===")
175     for r in metrics_rows:
176         print("{:<10s}  ARI={}  NMI={}  Sil={}".format(*r))
177     if best_by_ari[0] is not None:
178         print(f"\nBest by ARI: {best_by_ari[0]} (ARI={best_by_ari[1]:.4f})")
179
180     print(f"\nSaved figures & tables to: {outdir}")
181
182 def load_and_preprocess_brca(data_path: str, var_threshold: float = 1e-4,
183                             report_path: str = None):
184
185     # 1) Load with sample IDs in index
186     df = pd.read_csv(data_path, index_col=0)
187     n_samples_raw, n_cols_raw = df.shape
188     print(f"Raw data shape: {n_samples_raw} samples × {n_cols_raw} columns")
189
190     # 2) Normalize column names
191     df.columns = (
192         df.columns.astype(str)
193         .str.strip()
194         .str.replace("-", "_", regex=False)
195         .str.replace(" ", "_", regex=False)
196         .str.lower()
197     )

```



```

194 # 3) Detect clinical columns robustly
195 possible_clinical = ["subtype", "er_status", "pr_status", "her2_status", "
                        node", "metastasis"]
196 clinical_cols = [c for c in possible_clinical if c in df.columns]
197 if not clinical_cols:
198     raise ValueError("No clinical columns detected after normalization. "
199                      "Please check the CSV headers.")
200 print(f"Detected clinical columns: {clinical_cols}")
201
202 # 4) Split gene vs clinical
203 gene_cols = [c for c in df.columns if c not in clinical_cols]
204 # Force numeric on genes; coerce errors to NaN
205 df[gene_cols] = df[gene_cols].apply(pd.to_numeric, errors="coerce")
206
207 X_df = df[gene_cols].copy()
208 y_clinical = df[clinical_cols].copy()
209
210 # 5) Missing value stats and imputation
211 n_missing_total = int(X_df.isna().sum().sum())
212 n_rows_with_na = int(X_df.isna().any(axis=1).sum())
213 n_genes_with_na = int(X_df.isna().any(axis=0).sum())
214 print(f"Missing values in gene matrix: {n_missing_total} "
215       f"(rows with any NA: {n_rows_with_na}, genes with any NA: {
                n_genes_with_na})")
216
217 # Impute clinical: categorical by mode, numeric by mean
218 for col in y_clinical.columns:
219     if y_clinical[col].dtype == "O":
220         mode_vals = y_clinical[col].mode(dropna=True)
221         fill_val = mode_vals.iloc[0] if not mode_vals.empty else "Unknown"
222         y_clinical[col] = y_clinical[col].fillna(fill_val)
223     else:
224         y_clinical[col] = y_clinical[col].fillna(y_clinical[col].mean())
225
226 # Impute genes by column median (robust)
227 X_df = X_df.fillna(X_df.median())
228
229 # 6) Standardize (Z-score)
230 scaler = StandardScaler()
231 X_scaled = scaler.fit_transform(X_df)
232
233 # 7) Build and optionally save summary report
234 summary_lines = [
235     "=== BRCA Data Preprocessing Summary ===",
236     f"Original shape: {n_samples_raw} samples × {n_cols_raw} columns",
237     f"Detected clinical columns: {clinical_cols}",
238     f"Initial gene columns: {len(gene_cols)}",
239     f"Total missing values filled (genes): {n_missing_total}",
240     f"Rows with any NA (genes): {n_rows_with_na}",
241     f"Genes with any NA: {n_genes_with_na}",
242     f"Final standardized data shape: {X_scaled.shape}",
243 ]
244
245 # --- Clinical summary by type ---

```

```

246 summary_text = ["Clinical variable summary:"]
247 for col in y_clinical.columns:
248     if y_clinical[col].dtype == "O":
249         counts = y_clinical[col].value_counts(dropna=False)
250         summary_text.append(f"\n{col} (categorical):")
251         summary_text.append(str(counts))
252     else:
253         desc = y_clinical[col].describe()
254         summary_text.append(f"\n{col} (numeric):")
255         summary_text.append(str(desc))
256
257 report_text = "\n".join(summary_lines + ["\n".join(summary_text)])
258 print("\n" + report_text)
259
260 if report_path is not None:
261     report_dir = Path(report_path).parent
262     report_dir.mkdir(parents=True, exist_ok=True)
263     with open(report_path, "w", encoding="utf-8") as f:
264         f.write(report_text)
265     print(f"\nPreprocessing summary saved to {report_path}")
266
267 return X_scaled, y_clinical, gene_cols
268
269 def run_embedding(X, method='pca', n_components=2, seed=0, **kwargs):
270     method = method.lower()
271     if method == 'pca':
272         model = PCA(n_components=n_components, random_state=seed)
273         name = 'PCA'
274     elif method == 'nmf':
275         model = NMF(n_components=n_components, init='nndsvda', random_state=
276                     seed, max_iter=1000)
277         name = 'NMF'
278     elif method == 'tsne':
279         model = TSNE(n_components=n_components, random_state=seed, **kwargs)
280         name = 'tSNE'
281     elif method == 'umap':
282         model = umap.UMAP(n_components=n_components, random_state=seed, **
283                           kwargs)
284         name = 'UMAP'
285     elif method == 'spectral':
286         model = SpectralEmbedding(n_components=n_components, random_state=seed
287                                   )
288         name = 'Spectral'
289     else:
290         raise ValueError(f"Unknown method: {method}")
291
292     X_emb = model.fit_transform(X)
293     return X_emb, name
294
295 def _encode_categories(series):
296     cats = series.astype(str).fillna("NA").values
297     uniq = sorted(np.unique(cats))
298     mapping = {c: i for i, c in enumerate(uniq)}

```

```

297     idx = np.array([mapping[c] for c in cats], dtype=int)
298     return idx, mapping
299
300
301 def plot_embedding(X2d, meta, hue_col, title, outpath):
302     outpath = Path(outpath)
303     outpath.parent.mkdir(parents=True, exist_ok=True)
304
305     labels = meta[hue_col].astype(str).fillna("NA")
306     idx, mapping = _encode_categories(labels)
307     cmap = plt.cm.get_cmap('tab10', len(mapping))
308
309     plt.figure(figsize=(6.2, 5.2), dpi=120)
310     sc = plt.scatter(X2d[:,0], X2d[:,1], c=idx, s=14, cmap=cmap, alpha=0.85,
311                     edgecolors='none')
312     plt.title(f"{title} embedding colored by {hue_col}")
313     plt.xlabel("Dim 1"); plt.ylabel("Dim 2")
314     handles = [plt.Line2D([0],[0], marker='o', linestyle='',
315                           markersize=6, markerfacecolor=cmap(i),
316                           markeredgcolor='none')
317               for i in range(len(mapping))]
318     labels_sorted = list(mapping.keys())
319     plt.legend(handles, labels_sorted, bbox_to_anchor=(1.02, 1), loc='upper
320               left', fontsize=8, frameon=False)
321
322     plt.tight_layout()
323     plt.savefig(outpath, bbox_inches='tight')
324     plt.close()
325
326
327 def evaluate_embedding(X2d, subtype_series, seed=0):
328     mask = subtype_series.notna()
329     X_eval = X2d[mask.values]
330     y_ref = subtype_series[mask].astype(str).values
331
332     km = KMeans(n_clusters=5, random_state=seed, n_init=10)
333     pred = km.fit_predict(X_eval)
334
335     ari = adjusted_rand_score(y_ref, pred)
336     nmi = normalized_mutual_info_score(y_ref, pred)
337     sil = silhouette_score(X_eval, pred)
338
339     return ari, nmi, sil
340
341
342 if __name__ == "__main__":
343     import argparse
344     parser = argparse.ArgumentParser()
345     parser.add_argument("--data_path", type=str, default="Homework 1/Data/
346                       BRCA_data.csv")
347     parser.add_argument("--outdir", type=str, default="Homework 1/Latex/
348                       Results/Problem_2")
349     parser.add_argument("--seed", type=int, default=0)

```

```
345     args = parser.parse_args()
346
347     run_problem2_pipeline(
348         data_path=args.data_path,
349         outdir=args.outdir,
350         seed=args.seed
351     )
```