

STAT 5244 – Unsupervised Learning

Homework 2

Name: Chuyang Su UNI: cs4570

1 Mixture Models

1.1 EM Algorithm Derivation

Since we model count-valued data, assume each observation $x_i = (x_{i1}, \dots, x_{ip}) \in \mathbb{N}_0^p$ is generated from a finite mixture of *independent* Poisson distributions:

$$p(x_i; \pi, \lambda) = \sum_{k=1}^K \pi_k \prod_{j=1}^p \frac{e^{-\lambda_{kj}} \lambda_{kj}^{x_{ij}}}{x_{ij}!}, \quad \pi_k \geq 0, \quad \sum_{k=1}^K \pi_k = 1, \quad \lambda_{kj} > 0.$$

Here $\pi = (\pi_1, \dots, \pi_K)$ are mixture weights and $\lambda_k = (\lambda_{k1}, \dots, \lambda_{kp})$ are component-wise Poisson means.

Latent variables. Introduce latent indicators $z_{ik} \in \{0, 1\}$ with $\sum_{k=1}^K z_{ik} = 1$, where $z_{ik} = 1$ if x_i comes from component k . The complete-data likelihood is

$$L_c(\pi, \lambda) = \prod_{i=1}^n \prod_{k=1}^K \left[\pi_k \prod_{j=1}^p \frac{e^{-\lambda_{kj}} \lambda_{kj}^{x_{ij}}}{x_{ij}!} \right]^{z_{ik}}.$$

Taking logs and dropping constants independent of (π, λ) (i.e., $\log x_{ij}!$) gives the complete-data log-likelihood

$$\ell_c(\pi, \lambda) \propto \sum_{i=1}^n \sum_{k=1}^K z_{ik} \left[\log \pi_k + \sum_{j=1}^p (x_{ij} \log \lambda_{kj} - \lambda_{kj}) \right].$$

E-step Derivation. Since the latent indicators z_{ik} are unobserved, we take their conditional expectation under the current parameters. Define

$$\gamma_{ik} := \mathbb{E}[z_{ik} \mid x_i; \pi^{(t)}, \lambda^{(t)}] = P(z_{ik} = 1 \mid x_i; \pi^{(t)}, \lambda^{(t)}),$$

which represents the posterior probability that observation x_i belongs to component k .

Using Bayes' theorem,

$$P(z_{ik} = 1 \mid x_i; \pi^{(t)}, \lambda^{(t)}) = \frac{P(z_{ik} = 1; \pi^{(t)}) P(x_i \mid z_{ik} = 1; \lambda^{(t)})}{P(x_i; \pi^{(t)}, \lambda^{(t)})}.$$

Each term can be expressed as:

$$P(z_{ik} = 1; \pi^{(t)}) = \pi_k^{(t)}, \quad P(x_i \mid z_{ik} = 1; \lambda^{(t)}) = \prod_{j=1}^p \frac{e^{-\lambda_{kj}^{(t)}} (\lambda_{kj}^{(t)})^{x_{ij}}}{x_{ij}!},$$

and

$$P(x_i; \pi^{(t)}, \lambda^{(t)}) = \sum_{\ell=1}^K \pi_{\ell}^{(t)} \prod_{j=1}^p \frac{e^{-\lambda_{\ell j}^{(t)}} (\lambda_{\ell j}^{(t)})^{x_{ij}}}{x_{ij}!}.$$

Substituting these expressions into Bayes' rule yields:

$$\gamma_{ik} = \frac{\pi_k^{(t)} \prod_{j=1}^p e^{-\lambda_{kj}^{(t)}} (\lambda_{kj}^{(t)})^{x_{ij}} / x_{ij}!}{\sum_{\ell=1}^K \pi_{\ell}^{(t)} \prod_{j=1}^p e^{-\lambda_{\ell j}^{(t)}} (\lambda_{\ell j}^{(t)})^{x_{ij}} / x_{ij}!}.$$

Since the term $\prod_{j=1}^p x_{ij}!$ does not depend on k , it cancels out between numerator and denominator. Therefore, the final expression for the responsibilities is:

$$\gamma_{ik} = \frac{\pi_k^{(t)} \prod_{j=1}^p e^{-\lambda_{kj}^{(t)}} (\lambda_{kj}^{(t)})^{x_{ij}}}{\sum_{\ell=1}^K \pi_{\ell}^{(t)} \prod_{j=1}^p e^{-\lambda_{\ell j}^{(t)}} (\lambda_{\ell j}^{(t)})^{x_{ij}}}, \quad i = 1, \dots, n, \quad k = 1, \dots, K.$$

M-step. We maximize

$$Q(\pi, \lambda) = \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik} \left[\log \pi_k + \sum_{j=1}^p (x_{ij} \log \lambda_{kj} - \lambda_{kj}) \right]$$

subject to $\pi_k \geq 0$, $\sum_{k=1}^K \pi_k = 1$, and $\lambda_{kj} > 0$.

Update for π_k . Introduce a Lagrange multiplier η for the simplex constraint:

$$\mathcal{L}(\pi, \eta) = \sum_{k=1}^K \left(\sum_{i=1}^n \gamma_{ik} \right) \log \pi_k + \eta \left(1 - \sum_{k=1}^K \pi_k \right).$$

Setting the partial derivatives to zero,

$$\frac{\partial \mathcal{L}}{\partial \pi_k} = \frac{\sum_{i=1}^n \gamma_{ik}}{\pi_k} - \eta = 0 \quad \implies \quad \pi_k = \frac{\sum_{i=1}^n \gamma_{ik}}{\eta}.$$

Summing over k and using $\sum_{k=1}^K \pi_k = 1$ gives

$$1 = \sum_{k=1}^K \pi_k = \frac{1}{\eta} \sum_{k=1}^K \sum_{i=1}^n \gamma_{ik} = \frac{1}{\eta} \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik} = \frac{1}{\eta} \sum_{i=1}^n 1 = \frac{n}{\eta} \implies \eta = n.$$

Hence

$$\pi_k^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_{ik}.$$

(Concavity: $\partial^2 \mathcal{L} / \partial \pi_k^2 = -(\sum_i \gamma_{ik}) / \pi_k^2 < 0$.)

Update for λ_{kj} . For each (k, j) , the terms of Q that involve λ_{kj} are

$$Q_{kj}(\lambda_{kj}) = \sum_{i=1}^n \gamma_{ik} (x_{ij} \log \lambda_{kj} - \lambda_{kj}).$$

Differentiate and set to zero:

$$\frac{\partial Q_{kj}}{\partial \lambda_{kj}} = \sum_{i=1}^n \gamma_{ik} \left(\frac{x_{ij}}{\lambda_{kj}} - 1 \right) = 0 \quad \implies \quad \lambda_{kj} = \frac{\sum_{i=1}^n \gamma_{ik} x_{ij}}{\sum_{i=1}^n \gamma_{ik}}.$$

(Concavity: $\partial^2 Q_{kj} / \partial \lambda_{kj}^2 = -\sum_i \gamma_{ik} x_{ij} / \lambda_{kj}^2 < 0$ when $\lambda_{kj} > 0$.) Therefore

$$\lambda_{kj}^{(t+1)} = \frac{\sum_{i=1}^n \gamma_{ik} x_{ij}}{\sum_{i=1}^n \gamma_{ik}}.$$

1.2 Interpretation and Comparison of Poisson and Gaussian Mixture Models.

Both the Poisson Mixture Model (PMM) and the Gaussian Mixture Model (GMM) successfully converged and achieved almost identical clustering performance on the author–chapter word-frequency data.

The PMM converged in 16 iterations with a clustering purity of **0.8989**, while the GMM converged in 14 iterations with a purity of **0.9001**. The learned author–cluster mappings were consistent across models: Cluster 0 corresponds to *Shakespeare*, Cluster 1 to *London*, and Clusters 2–3 to *Austen*, indicating that Austen’s writing exhibits two stylistically distinct sub-modes.

Examining the estimated centroids revealed interpretable linguistic patterns. The Shakespeare cluster assigns high weights to archaic forms such as *thou*, *thy*, and *hath*, capturing the syntax of Early Modern English. London’s centroid emphasizes neutral verbs like *was*, *had*, and *were*, representing narrative realism, while Austen’s two clusters differ primarily in pronoun and modal-verb usage: one dominated by *she*, *her*, *would*, *could* (social dialogue tone), and the other by *my*, *our*, *only*, *such* (introspective narration).

By inspecting the soft responsibilities γ_{ik} , chapters with $\max_k \gamma_{ik} < 0.6$ were identified as *low-certainty chapters*. These typically occur near stylistic transitions or among authors with overlapping vocabularies. Such ambiguous sections highlight that soft clustering provides richer insights than hard assignments.

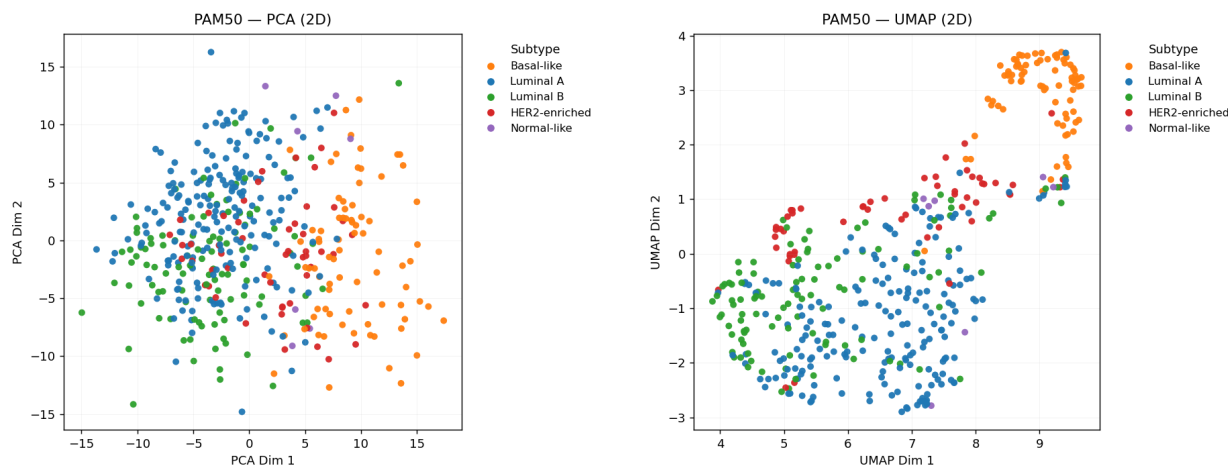
Overall, both mixture models achieve roughly 90% purity and uncover meaningful stylistic structures. While the GMM yields a marginally higher purity, the Poisson mixture remains theoretically more appropriate for discrete count data and offers comparable empirical performance.

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

2.1 Apply clustering techniques to explore this data set.

To explore the internal structure of the BRCA gene-expression dataset, I first applied a two-step dimensionality reduction approach using **PCA** followed by **UMAP**.

Based on conclusions drawn in Homework 1, this strategy provides an effective balance between global and local structure preservation: PCA removes redundant noise variables while retaining the main variance directions, and UMAP further compresses the PCA-transformed space, producing a low-dimensional embedding that reveals local neighborhood structure and facilitates visual interpretation. Following the empirical results from Homework 1, I selected 30 principal components (the “elbow” point in the cumulative variance curve) for PCA, and adopted UMAP parameters $n_neighbors = 10$, $min_dist = 0.0$, and $n_components = 2$ or 10, where the two-dimensional embedding was used for visualization and the ten-dimensional representation served as the feature space for clustering. The results are shown below:



(a) PAM50 subtypes visualized in PCA(2D) space.

(b) PAM50 subtypes visualized in UMAP(2D) space.

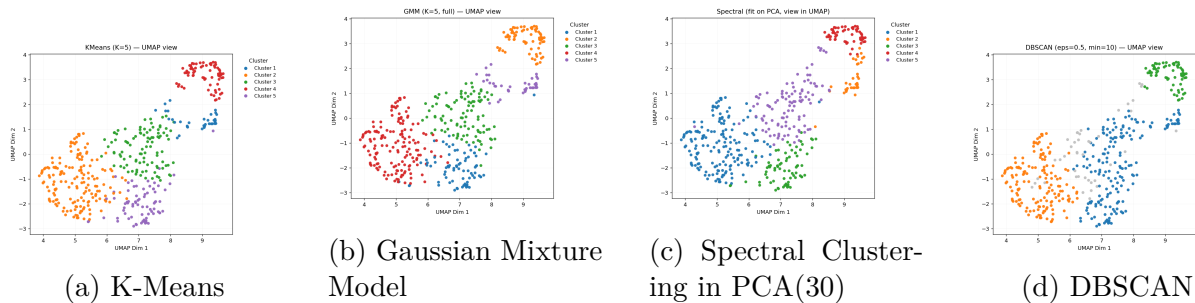
Subsequent clustering was performed using seven algorithms: **KMeans**, **Gaussian Mixture Model (GMM)**, **Spectral Clustering**, **DBSCAN**, and **Agglomerative Clustering** with three linkages (Ward, Average, and Complete). Given that the biological subtype reference (PAM50) contains five classes, all parametric clustering methods were set to $K = 5$, while DBSCAN automatically determined the number of clusters.

Among these, **Spectral Clustering** was treated as a special case. Because it internally constructs a similarity graph and performs Laplacian eigen-decomposition—conceptually similar to UMAP’s graph-based embedding—it is inappropriate to apply Spectral Clustering directly in UMAP space. Doing so would perform a second spectral decomposition on an already nonlinear manifold, likely distorting the intrinsic geometry. Therefore, Spectral Clustering was conducted on the PCA(30D) features instead, ensuring consistency with its

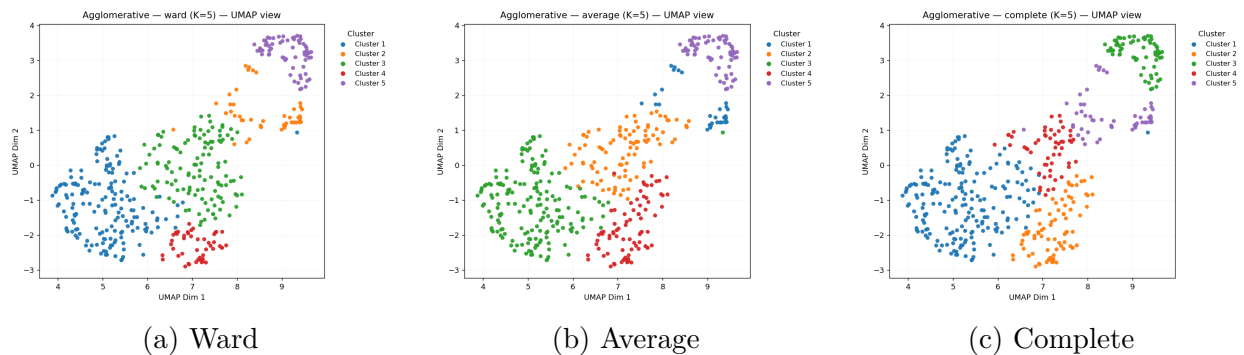
theoretical formulation.

For the Agglomerative Clustering, I compared only the **Ward**, **Average**, and **Complete** linkages, excluding **Single linkage** for three reasons: (1) it is highly sensitive to outliers and prone to chaining effects, which are amplified in high-dimensional continuous data; (2) its computational cost is high while its silhouette scores are unstable; and (3) the remaining three linkages already represent the major clustering behaviors.

All the clustering results expect of hierachical clusterings is shown below:



And three types of hierachical clusterings' results are:



Insights from the clustering results. Rather than directly comparing algorithmic performance, the focus here is on what each clustering method reveals about the *data manifold itself*. Across methods, two key structural patterns emerge. First, the dataset contains both compact and diffuse regions: **Basal-like** and **Normal-like** samples consistently form dense, isolated clusters, while **Luminal A** and **Luminal B** exhibit gradual overlap, suggesting a continuous transition rather than distinct boundaries. Second, the overall data distribution is non-spherical and heterogeneous in density, which explains why algorithms assuming isotropic clusters (KMeans, GMM, Ward) perform similarly and fail to separate overlapping subtypes.

DBSCAN highlights these density variations most clearly—it detects three main groups and labels the remaining sparse regions as noise. Although this underestimates the total number of subtypes, it accurately reflects the intrinsic density imbalance of the dataset: the rarest subtypes (**HER2-enriched**, **Normal-like**) are naturally absorbed as low-density regions. **Average linkage** produces smoother transitions that capture the gradual relationship between Luminal subtypes, while **Complete linkage** emphasizes inter-cluster separation. Taken together, these clustering outcomes indicate that the BRCA expression data contain

both discrete and continuous subtype structures—Basal-like being compact and distinct, whereas Luminal A/B lie on a continuum of transcriptional profiles—consistent with known biological heterogeneity among breast cancer subtypes.

2.2 Implement cluster validation techniques.

In this section, I evaluate the clustering results obtained in 2(a) through systematic parameter tuning and validation. For algorithms that require specifying the number of clusters K (*KMeans*, *GMM*, *Spectral Clustering*, and *Agglomerative Clustering* with three linkages), I performed a grid search over $K = 2, \dots, 9$ and computed three complementary validation criteria:

- **Silhouette Score** — measures intra-cluster cohesion and inter-cluster separation;
- **Stability** — estimated via bootstrap resampling, quantified by the median Adjusted Rand Index (ARI) between full-data and bootstrap partitions;
- **Generalizability** — assessed by a train–test split, where cluster assignments were propagated from training to testing samples (via centroids or k NN label transfer), and the Silhouette score on the test set was reported.

For **DBSCAN**, which is density-based and does not rely on K , a two-dimensional grid search was conducted across $\varepsilon \in \{0.3, 0.4, 0.5, 0.6, 0.7\}$ and $\text{min_samples} \in \{5, 10, 15\}$. Each configuration was evaluated using the same three validation criteria. Additionally, the k -distance graph was inspected to locate the elbow point corresponding to a suitable ε threshold.

All validation computations were performed on the same reduced feature spaces, PCA (30D) for Spectral Clustering, UMAP(10D) for others, to ensure consistency with the embeddings generated in Section 2(a).

Results of Validation Metrics

Table 1 summarizes the best-performing configurations for each clustering method according to the combined Silhouette, Stability, and Generalizability criteria.

Table 1: Best validation results per method.

Method	K	ε	min_samples	Silhouette	Stability	Generalizability
KMeans	2	–	–	0.549	0.994	0.468
GMM	2	–	–	0.551	0.983	0.473
Spectral (PCA)	2	–	–	0.205	0.959	0.217
Agglomerative (Ward)	2	–	–	0.558	0.673	0.481
Agglomerative (Average)	2	–	–	0.570	0.864	0.517
Agglomerative (Complete)	2	–	–	0.487	0.535	0.444
DBSCAN	–	0.3	15	0.292	0.929	0.344

Across all K -based algorithms, the optimal number of clusters was consistently $K = 2$. This indicates that the BRCA gene-expression manifold is dominated by two large-scale density regions rather than five fully separated clusters, which is biologically plausible since Luminal A/B and HER2-enriched subtypes form a continuous spectrum, while Basal-like and Normal-like remain distinct.

Discussion and Selection of the Best Result

Among all evaluated methods, the **Agglomerative Clustering (Average linkage, $K = 2$)** achieved the highest Silhouette (0.57) together with strong stability (0.86) and generalizability (0.52), providing the best trade-off between cluster compactness and robustness. KMeans and GMM yielded similar two-cluster partitions with excellent stability (> 0.98), indicating consistent global structure but limited ability to separate overlapping subtypes. Spectral Clustering performed worse due to over-smoothing in PCA space, and DBSCAN emphasized density variation rather than discrete boundaries, identifying a few compact regions while labeling sparse points as noise.

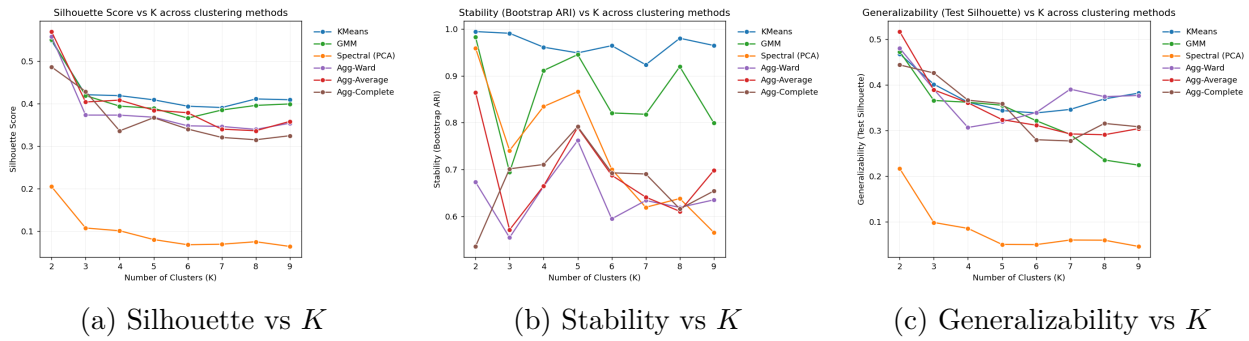


Figure 4: Validation metrics across different cluster numbers.

Figure 4 illustrates the variation of validation metrics with respect to the number of clusters K .

2.3 Interpret your cluster findings with respect to the metadata provided.

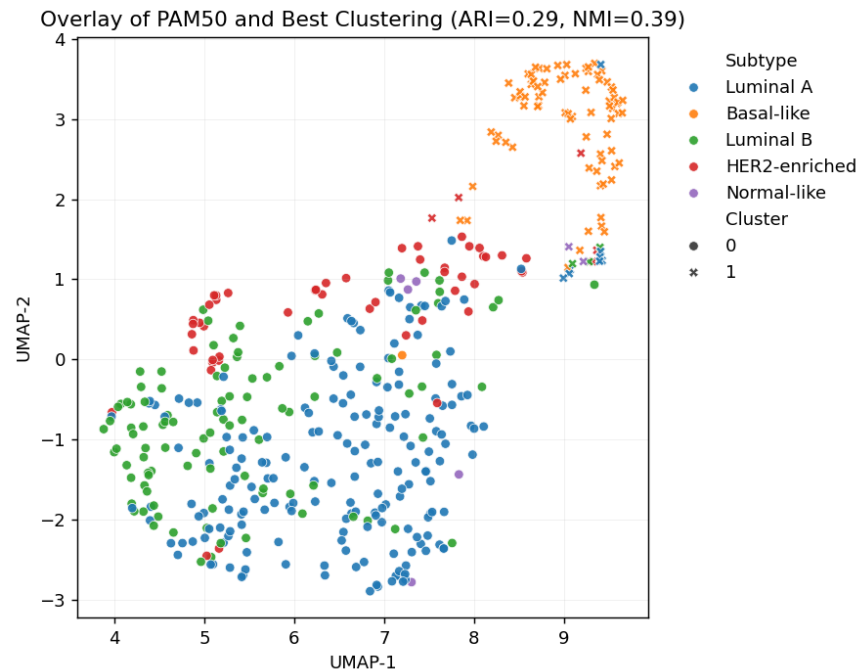


Figure 5: Overlay of PAM50 subtypes and the best clustering result (Agglomerative Average, $K = 2$) in UMAP(2D) space.

The two clusters discovered by the optimal Agglomerative (Average) model ($K = 2$) show a clear correspondence with the clinical PAM50 categories. As shown in Figure 5, one cluster primarily groups **Basal-like** and **Normal-like** tumors, while the other aggregates **Luminal A**, **Luminal B**, and **HER2-enriched** samples. This separation captures the well-known *luminal–basal dichotomy* in breast cancer biology: luminal tumors are typically hormone-receptor positive and exhibit similar expression patterns, whereas basal tumors represent the triple-negative subtype with distinct molecular characteristics. Therefore, the unsupervised clustering not only aligns with the known clinical taxonomy, but also highlights that the BRCA gene-expression landscape is dominated by two broad transcriptional regimes rather than five sharply separated classes.

A Appendix: Code Implementation

A.1 Problem 1b.

```

1  '''
2  Author: Chuyang Su cs4570@columbia.edu
3  Date: 2025-10-29 17:59:43
4  LastEditTime: 2025-10-30 14:06:41
5  FilePath: /Unsupervised-Learning-Homework/Homework 2/Code/Problem_1_b.py
6  Description:
7      EM algorithm for a mixture of Poisson distributions and fit to the author
8                                          data with K = 4 clusters.
9  '''
10 import os
11 import json
12 import numpy as np
13 import pandas as pd
14
15 DATA_PATH = r"Homework 2\Code\Data\authors.csv"
16 RESULT_DIR = r"Homework 2\Code\Result"
17 os.makedirs(RESULT_DIR, exist_ok=True)
18
19 def load_author_data(path):
20     df = pd.read_csv(path)
21     cols = list(df.columns)
22
23     author_col = "" if "" in cols else cols[0]
24
25     # Drop BookID column if exists
26     df = df.drop(columns=cols[-1])
27
28     y_names, y = np.unique(df[author_col].astype(str).values, return_inverse=
29                                     True)
30
31     feat_cols = [c for c in df.columns if c != author_col]
32     X = df[feat_cols].to_numpy(dtype=float)
33
34     return X, y, y_names, feat_cols
35
36 class PoissonMixtureResult:
37     __slots__ = ("pi", "lmbda", "gamma", "loglik_hist", "converged", "n_iter")
38     def __init__(self, pi, lmbda, gamma, loglik_hist, converged, n_iter):
39         self.pi = pi
40         self.lmbda = lmbda
41         self.gamma = gamma
42         self.loglik_hist = loglik_hist
43         self.converged = converged
44         self.n_iter = n_iter
45
46 def _softmax_logspace(logW, axis=1):
47     m = np.max(logW, axis=axis, keepdims=True)
48     W = np.exp(logW - m)
49     W /= W.sum(axis=axis, keepdims=True)
50     return W

```

```

48
49 def _logsumexp(A, axis=1):
50     m = np.max(A, axis=axis, keepdims=True)
51     return (m + np.log(np.sum(np.exp(A - m), axis=axis, keepdims=True))).
52                                     squeeze(axis)
53
54 def _init_params(X, K, random_state=None):
55     rng = np.random.default_rng(random_state)
56     n, p = X.shape
57     gamma = rng.dirichlet(alpha=np.ones(K), size=n)
58     Nk = gamma.sum(axis=0) + 1e-16
59     pi = Nk / Nk.sum()
60     lambda_ = (gamma.T @ X) / Nk[:, None]
61     lambda_ = np.clip(lambda_, 1e-8, None)
62     return pi, lambda_, gamma
63
64 def poisson_mixture_em(X, K, tol=1e-6, max_iter=500, random_state=None,
65                         verbose=False):
66     X = np.asarray(X, dtype=float)
67     if np.any(X < 0) or (np.abs(X - np.round(X)) > 1e-10).any():
68         raise ValueError("X must be nonnegative integer counts.")
69     n, p = X.shape
70     rng = np.random.default_rng(random_state)
71
72     pi, lambda_, gamma = _init_params(X, K, random_state=rng)
73     loglik_hist = []
74     eps = 1e-12
75
76     for it in range(1, max_iter + 1):
77         log_pi = np.log(np.clip(pi, eps, 1.0))
78         log_lambda = np.log(np.clip(lambda_, eps, None))
79         log_px_given_k = X @ log_lambda.T - np.sum(lambda_, axis=1)[None, :]
80         log_w = log_pi[None, :] + log_px_given_k
81         gamma = _softmax_logspace(log_w, axis=1)
82
83         Nk = gamma.sum(axis=0) + eps
84         pi = Nk / n
85         lambda_ = (gamma.T @ X) / Nk[:, None]
86         lambda_ = np.clip(lambda_, 1e-12, None)
87
88         ll_vec = _logsumexp(log_w, axis=1)
89         ll = float(ll_vec.sum())
90         loglik_hist.append(ll)
91
92         if it > 1:
93             ll_prev = loglik_hist[-2]
94             denom = max(1.0, abs(ll_prev))
95             if (ll - ll_prev) / denom < tol:
96                 return PoissonMixtureResult(pi, lambda_, gamma, loglik_hist,
97                                             True, it)
98
99     return PoissonMixtureResult(pi, lambda_, gamma, loglik_hist, False, max_iter)

```

```

98 def poisson_mixture_predict_proba(X, result):
99     X = np.asarray(X, dtype=float)
100     log_pi = np.log(np.clip(result.pi, 1e-12, 1.0))
101     log_lambda = np.log(np.clip(result.lmbda, 1e-12, None))
102     log_px_given_k = X @ log_lambda.T - np.sum(result.lmbda, axis=1)[None, :]
103     log_w = log_pi[None, :] + log_px_given_k
104     return _softmax_logspace(log_w, axis=1)
105
106 def poisson_mixture_predict(X, result):
107     return np.argmax(poisson_mixture_predict_proba(X, result), axis=1)
108
109 # Main function
110 if __name__ == "__main__":
111     X, y, y_names, feat_cols = load_author_data(DATA_PATH)
112
113     K = len(y_names)
114     res = poisson_mixture_em(X, K=K, tol=1e-6, max_iter=1000, random_state=0,
115                             verbose=False)
116
117     hard = poisson_mixture_predict(X, res)
118
119     cont = np.zeros((K, K), dtype=int)
120     for yi, hi in zip(y, hard):
121         cont[yi, hi] += 1
122
123     purity = float(cont.max(axis=0).sum() / len(y))
124     cluster_to_author = {int(k): str(y_names[int(np.argmax(cont[:, k]))]) for
125                          k in range(K)}
126
127     threshold = 0.6
128     max_resp = res.gamma.max(axis=1)
129     low_certainty_idx = np.where(max_resp < threshold)[0].tolist()
130
131     low_certainty_info = []
132     for i in low_certainty_idx:
133         low_certainty_info.append({
134             "chapter_index": int(i),
135             "author_true": str(y_names[y[i]]),
136             "pred_cluster": int(hard[i]),
137             "max_gamma": float(max_resp[i]),
138             "responsibilities": [float(v) for v in res.gamma[i]]
139         })
140
141     out_path = os.path.join(RESULT_DIR, f"poisson_mixture_result.json")
142     result_json = {
143         "converged": bool(res.converged),
144         "n_iter": int(res.n_iter),
145         "final_loglik": float(res.loglik_hist[-1]),
146         "n_authors": int(K),
147         "authors": y_names.tolist(),
148         "feature_columns": feat_cols,
149         "em": {
150             "pi": res.pi.astype(float).tolist(),
151             "lambda": res.lmbda.astype(float).tolist(),
152             "loglik_hist": [float(v) for v in res.loglik_hist]
153         },

```

```

150     "predictions": {
151         "cluster_pred": [int(v) for v in hard],
152         "responsibilities": res.gamma.astype(float).tolist()
153     },
154     "validation": {
155         "purity": float(purity),
156         "contingency": cont.astype(int).tolist(),
157         "cluster_sizes": np.bincount(hard, minlength=K).astype(int).tolist(),
158         "cluster_to_author_map": {str(k): v for k, v in cluster_to_author.items()},
159         "low_certainty_chapters": low_certainty_info,
160         "threshold": float(threshold)
161     }
162 }
163
164 with open(out_path, "w", encoding="utf-8") as f:
165     json.dump(result_json, f, ensure_ascii=False, indent=2)

```

A.2 Problem 1c.

```

1  '''
2  Author: Chuyang Su cs4570@columbia.edu
3  Date: 2025-10-29 21:47:48
4  LastEditTime: 2025-10-30 14:15:24
5  FilePath: /Unsupervised-Learning-Homework/Homework 2/Code/Problem_1_c.py
6  Description:
7      Fit a Gaussian mixture model to the author data.
8  '''
9  import os
10 import json
11 import numpy as np
12 import pandas as pd
13
14 DATA_PATH = r"Homework 2\Code\Data\authors.csv"
15 RESULT_DIR = r"Homework 2\Code\Result"
16 os.makedirs(RESULT_DIR, exist_ok=True)
17
18 from Problem_1_b import load_author_data
19
20 class GMMResult:
21     __slots__ = ("pi", "mu", "var", "gamma", "loglik_hist", "converged", "n_iter")
22
23     def __init__(self, pi, mu, var, gamma, loglik_hist, converged, n_iter):
24         self.pi = pi          # (K,)
25         self.mu = mu          # (K, p)
26         self.var = var        # (K, p) diagonal covariances
27         self.gamma = gamma    # (n, K)
28         self.loglik_hist = loglik_hist
29         self.converged = converged
30         self.n_iter = n_iter

```

```

31 def _log_gauss_diag(X, mu, var):
32     eps = 1e-10
33     var = np.clip(var, eps, None)
34     n, p = X.shape
35     K = mu.shape[0]
36     log_det = np.sum(np.log(var), axis=1)
37     inv_var = 1.0 / var
38     quad = ((X[:, None, :] - mu[None, :, :]) ** 2 * inv_var[None, :, :]).sum(
39         axis=2)
40     return -0.5 * (p * np.log(2.0 * np.pi) + log_det[None, :] + quad)
41
42 def _softmax_logspace(logW, axis=1):
43     m = np.max(logW, axis=axis, keepdims=True)
44     W = np.exp(logW - m)
45     W /= W.sum(axis=axis, keepdims=True)
46     return W
47
48 def _logsumexp(A, axis=1):
49     m = np.max(A, axis=axis, keepdims=True)
50     return (m + np.log(np.sum(np.exp(A - m), axis=axis, keepdims=True))).squeeze(axis)
51
52 def _init_gmm_params(X, K, random_state=None):
53     rng = np.random.default_rng(random_state)
54     n, p = X.shape
55     mu = np.empty((K, p))
56     idx0 = rng.integers(0, n)
57     mu[0] = X[idx0]
58     d2 = np.sum((X - mu[0])**2, axis=1) + 1e-12
59     for k in range(1, K):
60         probs = d2 / d2.sum()
61         idx = rng.choice(n, p=probs)
62         mu[k] = X[idx]
63         d2 = np.minimum(d2, np.sum((X - mu[k])**2, axis=1) + 1e-12)
64     dist2 = ((X[:, None, :] - mu[None, :, :])**2).sum(axis=2)
65     gamma = np.zeros((n, K))
66     gamma[np.arange(n), np.argmin(dist2, axis=1)] = 1.0
67     Nk = gamma.sum(axis=0) + 1e-12
68     pi = Nk / n
69     var = (gamma.T @ (X**2)) / Nk[:, None] - ((gamma.T @ X) / Nk[:, None])**2
70     var = np.clip(var, 1e-6, None)
71     return pi, mu, var, gamma
72
73 def gaussian_mixture_em(
74     X, K, tol=1e-6, max_iter=500, random_state=None, verbose=False
75 ):
76     X = np.asarray(X, dtype=float)
77     n, p = X.shape
78     eps = 1e-12
79     pi, mu, var, gamma = _init_gmm_params(X, K, random_state=random_state)
80     loglik_hist = []
81

```

```

82     for it in range(1, max_iter + 1):
83         # E-step
84         log_pi = np.log(np.clip(pi, eps, 1.0)) # (K,)
85         log_px_given_k = _log_gauss_diag(X, mu, var) # (n,K)
86         log_w = log_pi[None, :] + log_px_given_k
87         gamma = _softmax_logspace(log_w, axis=1) # (n,K)
88
89         # M-step
90         Nk = gamma.sum(axis=0) + eps # (K,)
91         pi = Nk / n
92         mu = (gamma.T @ X) / Nk[:, None] # (K,p)
93         # diagonal covariance
94         Ex2 = (gamma.T @ (X**2)) / Nk[:, None] # (K,p)
95         var = Ex2 - mu**2
96         var = np.clip(var, 1e-6, None)
97
98         # Log-likelihood
99         ll = float(_logsumexp(log_w, axis=1).sum())
100        loglik_hist.append(ll)
101
102        if verbose and (it == 1 or it % 10 == 0):
103            print(f"Iter {it:4d} loglik={ll:.6f}")
104
105        if it > 1:
106            ll_prev = loglik_hist[-2]
107            denom = max(1.0, abs(ll_prev))
108            if (ll - ll_prev) / denom < tol:
109                if verbose:
110                    print(f"Converged at iter {it} Δ rel={(ll-ll_prev)/denom:.3e}")
111                return GMMResult(pi, mu, var, gamma, loglik_hist, True, it)
112
113        if verbose:
114            print("Reached max_iter without convergence.")
115        return GMMResult(pi, mu, var, gamma, loglik_hist, False, max_iter)
116
117    def gmm_predict_proba(X, result):
118        log_pi = np.log(np.clip(result.pi, 1e-12, 1.0))
119        log_px_given_k = _log_gauss_diag(X, result.mu, result.var)
120        log_w = log_pi[None, :] + log_px_given_k
121        return _softmax_logspace(log_w, axis=1)
122
123    def gmm_predict(X, result):
124        return np.argmax(gmm_predict_proba(X, result), axis=1)
125
126    # Main function
127    if __name__ == "__main__":
128        X, y, y_names, feat_cols = load_author_data(DATA_PATH)
129        K = len(y_names)
130
131        res = gaussian_mixture_em(
132            X, K=K, tol=1e-6, max_iter=1000, random_state=0, verbose=False
133        )
134        hard = gmm_predict(X, res)

```

```

135
136     cont = np.zeros((K, K), dtype=int)
137     for yi, hi in zip(y, hard):
138         cont[yi, hi] += 1
139     purity = float(cont.max(axis=0).sum() / len(y))
140     cluster_to_author = {
141         str(k): str(y_names[int(np.argmax(cont[:, k]))]) for k in range(K)
142     }
143
144     threshold = 0.6
145     max_resp = res.gamma.max(axis=1)
146     low_idx = np.where(max_resp < threshold)[0].tolist()
147     low_certainty_info = []
148     for i in low_idx:
149         low_certainty_info.append({
150             "chapter_index": int(i),
151             "author_true": str(y_names[y[i]]),
152             "pred_cluster": int(hard[i]),
153             "max_gamma": float(max_resp[i]),
154             "responsibilities": [float(v) for v in res.gamma[i]]
155         })
156
157     out_path = os.path.join(RESULT_DIR, f"gmm_result.json")
158
159     result_json = {
160         "converged": bool(res.converged),
161         "n_iter": int(res.n_iter),
162         "final_loglik": float(res.loglik_hist[-1]),
163         "n_authors": int(K),
164         "authors": y_names.tolist(),
165         "feature_columns": feat_cols,
166         "em": {
167             "pi": res.pi.astype(float).tolist(),
168             "mu": res.mu.astype(float).tolist(),
169             "var_diag": res.var.astype(float).tolist(),
170             "loglik_hist": [float(v) for v in res.loglik_hist]
171         },
172         "predictions": {
173             "cluster_pred": [int(v) for v in hard],
174             "responsibilities": res.gamma.astype(float).tolist()
175         },
176         "validation": {
177             "purity": float(purity),
178             "contingency": cont.astype(int).tolist(),
179             "cluster_sizes": np.bincount(hard, minlength=K).astype(int).tolist(),
180             "cluster_to_author_map": {str(k): v for k, v in cluster_to_author.items()},
181             "low_certainty_chapters": low_certainty_info,
182             "threshold": float(threshold)
183         }
184     }
185
186     with open(out_path, "w", encoding="utf-8") as f:

```

```
187 json.dump(result_json, f, ensure_ascii=False, indent=2)
```

A.3 Problem 2a.

```
1 '''
2 Author: Chuyang Su cs4570@columbia.edu
3 Date: 2025-10-30 14:36:11
4 LastEditTime: 2025-10-31 02:14:40
5 FilePath: /Unsupervised-Learning-Homework/Homework 2/Code/Problem_2_a.py
6 Description:
7     Apply clustering techniques(KMeans, GMM, Spectral Clustering,
8                                     Agglomerative Clustering, DBSCAN)
9                                     to explore the BRCA gene expression
10                                    data set.
11 '''
12 import os
13 os.environ.setdefault("OMP_NUM_THREADS", "2")
14 os.environ.setdefault("MKL_NUM_THREADS", "2")
15 os.environ.setdefault("LOKY_MAX_CPU_COUNT", "8")
16
17 import warnings
18 from matplotlib import MatplotlibDeprecationWarning
19 warnings.filterwarnings("ignore", category=UserWarning,
20                         message="n_jobs value 1 overridden to 1 by setting
21                                 random_state",
22                                 module=r"umap\.umap_")
23 warnings.filterwarnings("ignore", category=UserWarning,
24                         message="KMeans is known to have a memory leak on
25                                 Windows with
26                                 MKL",
27                                 module=r"sklearn\.cluster\._kmeans")
28 warnings.filterwarnings("ignore", category=UserWarning,
29                         message="Could not find the number of physical cores",
30                         module=r"joblib\.externals\.loky\.backend\.context")
31 warnings.filterwarnings("ignore", category=MatplotlibDeprecationWarning)
32
33 from pathlib import Path
34 import numpy as np
35 import pandas as pd
36 import matplotlib.pyplot as plt
37 import matplotlib as mpl
38
39 from sklearn.preprocessing import StandardScaler
40 from sklearn.decomposition import PCA
41 from sklearn.cluster import KMeans, SpectralClustering,
42                               AgglomerativeClustering, DBSCAN
43 from sklearn.mixture import GaussianMixture
44 from packaging.version import parse as vparse
45 import sklearn
46 import umap
47
48 def load_brca(data_path: str):
```



```

42 df = pd.read_csv(data_path, index_col=0)
43 df.columns = (
44     df.columns.astype(str)
45     .str.strip()
46     .str.replace("-", "_", regex=False)
47     .str.replace(" ", "_", regex=False)
48     .str.lower()
49 )
50 clinical_cols = [c for c in ["subtype", "er_status", "pr_status", "
                             her2_status", "node", "metastasis"]
51                  if c in df.columns]
52 gene_cols = [c for c in df.columns if c not in clinical_cols]
53
54 X = df[gene_cols].apply(pd.to_numeric, errors="coerce")
55 X = X.fillna(X.median())
56 X = StandardScaler().fit_transform(X.values)
57
58 meta = df[clinical_cols].copy()
59 return X, meta
60
61 def _encode(series):
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 PAM50_COLOR_MAP = {
69     "Luminal A": "#1f77b4",      # blue
70     "Luminal B": "#2ca02c",      # green
71     "Basal-like": "#ff7f0e",     # orange
72     "HER2-enriched": "#d62728",  # red
73     "Normal-like": "#9467bd",    # purple
74 }
75
76 def scatter_by_labels(
77     X2,
78     labels,
79     title,
80     outpath,
81     xlab="Dim 1",
82     ylab="Dim 2",
83     color_map=None,
84     order=None,
85     s=28,
86     legend_title=None,
87 ):
88     from pathlib import Path
89     outpath = Path(outpath); outpath.parent.mkdir(parents=True, exist_ok=True)
90
91     lbl = pd.Series(labels)
92     if pd.api.types.is_numeric_dtype(lbl):
93         uniq_vals = sorted(lbl.unique())
94         if -1 in uniq_vals:

```

```

95     uniq_vals = [v for v in uniq_vals if v != -1] + [-1]
96     uniq = [str(v) for v in uniq_vals]
97     lbl_str = lbl.astype(str)
98     legend_names = []
99     for u in uniq_vals:
100         if u == -1:
101             legend_names.append("Noise")
102         else:
103             legend_names.append(f"Cluster {int(u)+1}")
104     else:
105         lbl_str = lbl.astype(str).fillna("NA")
106         if order is None:
107             uniq = sorted(lbl_str.unique())
108         else:
109             present = [u for u in order if str(u) in set(lbl_str)]
110             rest = [u for u in lbl_str.unique() if u not in present]
111             uniq = present + sorted(rest)
112         legend_names = uniq
113
114     colors = []
115     if color_map is not None:
116         for u in uniq:
117             key = u if u in color_map else u.lower()
118             colors.append(color_map.get(key, "#9E9E9E"))
119     else:
120         # clean, professional palette (Tableau 20-ish extended) + reserved
121         #                                     gray for Noise
122         base = ["#1f77b4", "#ff7f0e", "#2ca02c", "#d62728", "#9467bd",
123               "#8c564b", "#e377c2", "#7f7f7f", "#bcbd22", "#17becf",
124               "#4c78a8", "#f58518", "#54a24b", "#e45756", "#72b7b2",
125               "#eecc3b", "#b279a2", "#ff9da6", "#9d755d", "#bab0ac"]
126         # if there is a "Noise" class, make it light gray
127         if pd.api.types.is_numeric_dtype(lbl) and (-1 in lbl.unique()):
128             # map last label (Noise) to gray
129             need = len(uniq) - 1
130             while len(base) < need:
131                 base = base + base
132             colors = base[:need] + ["#C0C0C0"]
133         else:
134             need = len(uniq)
135             while len(base) < need:
136                 base = base + base
137             colors = base[:need]
138
139         # map label string -> palette index
140         idx_map = {u: i for i, u in enumerate(uniq)}
141         idx = lbl_str.map(idx_map).values
142         point_colors = [colors[i] for i in idx]
143
144     plt.figure(figsize=(7.2, 5.6), dpi=130)
145     plt.scatter(X2[:, 0], X2[:, 1], c=point_colors, s=s, alpha=0.95,
146               edgecolors='none')
147     plt.xlabel(xlab); plt.ylabel(ylab); plt.title(title, pad=8)
148     plt.grid(True, linewidth=0.3, alpha=0.25)

```

```

147
148 handles = [plt.Line2D([0], [0], marker='o', linestyle='',
149                       markersize=6, markerfacecolor=colors[i],
150                       markeredgecolor='none') for i in range(len(uniq))]
151 if legend_title is None:
152     legend_title = "Cluster" if pd.api.types.is_numeric_dtype(lbl) else "
                                Label"
153 if len(legend_names) <= 25:
154     plt.legend(handles, legend_names,
155               bbox_to_anchor=(1.02, 1), loc='upper left',
156               fontsize=9, frameon=False, title=legend_title)
157
158 plt.tight_layout()
159 plt.savefig(outpath, bbox_inches='tight')
160 plt.close()
161
162 def plot_pam50(X2, pam50_series, title, outpath, xlab, ylab):
163     outpath.parent.mkdir(parents=True, exist_ok=True)
164
165     # enforce canonical label order
166     ordered_labels = ["Basal-like", "Luminal A", "Luminal B", "HER2-enriched",
                                "Normal-like"]
167     pam50 = pam50_series.astype(str).fillna("NA")
168
169     # color list in that order
170     colors = [PAM50_COLOR_MAP.get(lbl, "#999999") for lbl in ordered_labels]
171
172     # map label -> index for plotting
173     idx_map = {lbl: i for i, lbl in enumerate(ordered_labels)}
174     idx = [idx_map.get(v, -1) for v in pam50]
175     point_colors = [colors[i] if i >= 0 else "#999999" for i in idx]
176
177     plt.figure(figsize=(7, 5.5), dpi=130)
178     plt.scatter(X2[:, 0], X2[:, 1], c=point_colors, s=30, alpha=0.95,
                                edgecolors="none")
179     plt.xlabel(xlab); plt.ylabel(ylab)
180     plt.title(title, pad=8)
181     plt.grid(True, linewidth=0.3, alpha=0.25)
182
183     # legend - fixed text order and color
184     handles = [plt.Line2D([0], [0], marker='o', linestyle='',
185                           markersize=6, markerfacecolor=colors[i],
186                           markeredgecolor='none')
187               for i in range(len(ordered_labels))]
188     legend_labels = ["Basal-like", "Luminal A", "Luminal B",
189                     "HER2-enriched", "Normal-like"]
190     plt.legend(handles, legend_labels,
191               bbox_to_anchor=(1.02, 1), loc='upper left',
192               fontsize=9, frameon=False, title="Subtype")
193
194     plt.tight_layout()
195     plt.savefig(outpath, bbox_inches="tight")
196     plt.close()
197

```

```

198 def compute_embeddings(X, seed=0, n_pcs_feat=30):
199     # PCA for visualization (2D) and as preprocessor for UMAP
200     X_pca2 = PCA(n_components=2, random_state=seed).fit_transform(X)
201     X_pca30 = PCA(n_components=min(n_pcs_feat, X.shape[1]), random_state=seed)
202                 .fit_transform(X)
203
204     # UMAP on PCA-30
205     u10 = umap.UMAP(n_components=10, n_neighbors=10, min_dist=0.0,
206                     random_state=seed, metric="euclidean")
207     X_umap10 = u10.fit_transform(X_pca30)
208
209     u2 = umap.UMAP(n_components=2, n_neighbors=10, min_dist=0.0,
210                   random_state=seed, metric="euclidean")
211     X_umap2 = u2.fit_transform(X_pca30)
212
213     return X_pca2, X_pca30, X_umap10, X_umap2
214
215 def clustering_all(X_umap10, X_umap2, X_spectral_input, outdir: Path, seed=0):
216     k = 5 # PAM50 count
217
218     labels = KMeans(n_clusters=k, n_init=20, random_state=seed).fit_predict(
219                     X_umap10)
220
221     scatter_by_labels(
222         X_umap2, labels,
223         "KMeans (K=5) - UMAP view",
224         outdir / "cluster__kmeans_k5_umap2.png",
225         "UMAP Dim 1", "UMAP Dim 2",
226         legend_title="Cluster"
227     )
228
229     # Spectral on PCA
230     labels = SpectralClustering(
231         n_clusters=k, affinity="nearest_neighbors",
232         n_neighbors=10, assign_labels="kmeans",
233         random_state=seed, n_init=10
234     ).fit_predict(X_spectral_input)
235
236     scatter_by_labels(
237         X_umap2, labels,
238         "Spectral (fit on PCA, view in UMAP)",
239         outdir / "cluster__spectral_pca_umap2.png",
240         "UMAP Dim 1", "UMAP Dim 2",
241         legend_title="Cluster"
242     )
243
244     labels = GaussianMixture(
245         n_components=k, covariance_type="full",
246         random_state=seed, n_init=5
247     ).fit_predict(X_umap10)
248
249     scatter_by_labels(
250         X_umap2, labels,
251         "GMM (K=5, full) - UMAP view",
252         outdir / "cluster__gmm_k5_umap2.png",
253         "UMAP Dim 1", "UMAP Dim 2",
254         legend_title="Cluster"

```

```

250 )
251
252 labels = DBSCAN(eps=0.5, min_samples=10, metric="euclidean").fit_predict(
253                                     X_umap10)
254
255 scatter_by_labels(
256     X_umap2, labels,
257     "DBSCAN (eps=0.5, min=10) - UMAP view",
258     outdir / "cluster__dbscan_eps0p5_min10_umap2.png",
259     "UMAP Dim 1", "UMAP Dim 2",
260     legend_title="Cluster"
261 )
262
263 new_api = vparse(sklearn.__version__) >= vparse("1.2")
264
265 labels = AgglomerativeClustering(n_clusters=k, linkage="ward").fit_predict(
266                                     X_umap10)
267
268 scatter_by_labels(
269     X_umap2, labels,
270     "Agglomerative - ward (K=5) - UMAP view",
271     outdir / "cluster__agg_ward_k5_umap2.png",
272     "UMAP Dim 1", "UMAP Dim 2",
273     legend_title="Cluster"
274 )
275
276 if new_api:
277     model = AgglomerativeClustering(n_clusters=k, linkage="average",
278                                     metric="euclidean")
279
280 else:
281     model = AgglomerativeClustering(n_clusters=k, linkage="average",
282                                     affinity="euclidean")
283
284 labels = model.fit_predict(X_umap10)
285
286 scatter_by_labels(
287     X_umap2, labels,
288     "Agglomerative - average (K=5) - UMAP view",
289     outdir / "cluster__agg_average_k5_umap2.png",
290     "UMAP Dim 1", "UMAP Dim 2",
291     legend_title="Cluster"
292 )
293
294 if new_api:
295     model = AgglomerativeClustering(n_clusters=k, linkage="complete",
296                                     metric="euclidean")
297
298 else:
299     model = AgglomerativeClustering(n_clusters=k, linkage="complete",
300                                     affinity="euclidean")
301
302 labels = model.fit_predict(X_umap10)
303
304 scatter_by_labels(
305     X_umap2, labels,
306     "Agglomerative - complete (K=5) - UMAP view",
307     outdir / "cluster__agg_complete_k5_umap2.png",
308     "UMAP Dim 1", "UMAP Dim 2",
309     legend_title="Cluster"
310 )

```

```

298
299 def run(
300     data_path="Homework 2/Code/Data/BRCA_data.csv",
301     outdir="Homework 2/Code/Result/Problem_2",
302     seed=25
303 ):
304     outdir = Path(outdir); outdir.mkdir(parents=True, exist_ok=True)
305
306     # data
307     X, meta = load_brca(data_path)
308     pam50 = meta["subtype"] if "subtype" in meta.columns else pd.Series(["NA"]
309                                     * X.shape[0])
310
311     # embeddings (PCA-2 and UMAP-2 are the only two embedding figures we save)
312     X_pca2, X_pca30, X_umap10, X_umap2 = compute_embeddings(X, seed=seed,
313                                     n_pcs_feat=30)
314     plot_pam50(X_pca2, pam50, "PAM50 - PCA (2D)",
315                 outdir / "pam50__pca2.png", "PCA Dim 1", "PCA Dim 2")
316     plot_pam50(X_umap2, pam50, "PAM50 - UMAP (2D)",
317                 outdir / "pam50__umap2.png", "UMAP Dim 1", "UMAP Dim 2")
318     clustering_all(X_umap10, X_umap2, X_pca30, outdir, seed=seed)
319
320 if __name__ == "__main__":
321     import argparse
322     parser = argparse.ArgumentParser()
323     parser.add_argument("--data_path", type=str, default="Homework 2/Code/Data
324                                     /BRCA_data.csv")
325     parser.add_argument("--outdir", type=str, default="Homework 2\Latex\
326                                     Figures")
327     parser.add_argument("--seed", type=int, default=25)
328     args = parser.parse_args()
329     run(data_path=args.data_path, outdir=args.outdir, seed=args.seed)

```

A.4 Problem 2b

```

1  '''
2  Author: Chuyang Su cs4570@columbia.edu
3  Date: 2025-10-31 02:09:44
4  LastEditTime: 2025-10-31 02:50:01
5  FilePath: /Unsupervised-Learning-Homework/Homework 2/Code/Problem_2_b.py
6  Description:
7      Validation.
8  '''
9  import os
10 os.environ.setdefault("OMP_NUM_THREADS", "2")
11 os.environ.setdefault("MKL_NUM_THREADS", "2")
12 os.environ.setdefault("LOKY_MAX_CPU_COUNT", "8")
13
14 import warnings
15 warnings.filterwarnings("ignore")
16

```

```

17 from pathlib import Path
18 import numpy as np
19 import pandas as pd
20 import matplotlib.pyplot as plt
21 import seaborn as sns
22 from pathlib import Path
23
24 from packaging.version import parse as vparse
25 import sklearn
26 from sklearn.cluster import KMeans, SpectralClustering,
27                               AgglomerativeClustering, DBSCAN
28 from sklearn.mixture import GaussianMixture
29 from sklearn.metrics import silhouette_score, adjusted_rand_score
30 from sklearn.neighbors import KNeighborsClassifier
31 from sklearn.utils import check_random_state
32
33 from Problem_2_a import load_brca, compute_embeddings
34
35 def _kmeans_fit_predict(X, k, seed):
36     return KMeans(n_clusters=k, n_init=20, random_state=seed).fit_predict(X)
37
38 def _gmm_fit_predict(X, k, seed):
39     gm = GaussianMixture(n_components=k, covariance_type="full", random_state=
40                           seed, n_init=5)
41     return gm.fit_predict(X), gm
42
43 def _spectral_fit_predict(X_pca30, k, seed):
44     return SpectralClustering(
45         n_clusters=k, affinity="nearest_neighbors",
46         n_neighbors=10, assign_labels="kmeans",
47         random_state=seed, n_init=10
48     ).fit_predict(X_pca30)
49
50 def _agglo_fit_predict(X, k, linkage):
51     new_api = vparse(sklearn.__version__) >= vparse("1.2")
52     if linkage == "ward":
53         model = AgglomerativeClustering(n_clusters=k, linkage="ward")
54     else:
55         if new_api:
56             model = AgglomerativeClustering(n_clusters=k, linkage=linkage,
57                                             metric="euclidean")
58         else:
59             model = AgglomerativeClustering(n_clusters=k, linkage=linkage,
60                                             affinity="euclidean")
61     return model.fit_predict(X)
62
63 def _dbscan_fit_predict(X, eps, min_samples):
64     return DBSCAN(eps=eps, min_samples=min_samples, metric="euclidean").
65         fit_predict(X)
66
67 def _silhouette_safe(X, labels):
68     labs = np.asarray(labels)
69     if len(np.unique(labs)) <= 1 or np.all(labs == -1):
70         return np.nan

```

```

66     try:
67         return silhouette_score(X, labs)
68     except Exception:
69         return np.nan
70
71 def _bootstrap_stability(X, fit_fn, n_boot=10, seed=0):
72     rng = check_random_state(seed)
73     full_labels = fit_fn(X)
74     if len(np.unique(full_labels)) <= 1:
75         return np.nan
76     n = X.shape[0]
77     aris = []
78     for _ in range(n_boot):
79         idx = rng.choice(n, size=int(0.8 * n), replace=False)
80         boot_labels = fit_fn(X[idx])
81         if len(np.unique(boot_labels)) <= 1:
82             aris.append(np.nan)
83         else:
84             aris.append(adjusted_rand_score(full_labels[idx], boot_labels))
85     aris = np.array(aris, dtype=float)
86     return float(np.nanmedian(aris))
87
88 def _generalizability_test_silhouette(X, labels, tag, k=None, seed=0):
89     """
90     80/20 split: fit clusterer on train, assign test labels via:
91     - kmeans/gmm: model predict
92     - spectral/agg/dbscan: kNN(label transfer) from train to test
93     Return test silhouette.
94     """
95     rng = check_random_state(seed)
96     n = X.shape[0]
97     perm = rng.permutation(n)
98     cut = int(0.8 * n)
99     tr, te = perm[:cut], perm[cut:]
100    Xtr, Xte = X[tr], X[te]
101
102    # fit on train
103    if tag == "kmeans":
104        km = KMeans(n_clusters=k, n_init=20, random_state=seed).fit(Xtr)
105        yte = km.predict(Xte)
106    elif tag == "gmm":
107        gm = GaussianMixture(n_components=k, covariance_type="full",
108                             random_state=seed, n_init=5).
109                             fit(Xtr)
110
111        yte = gm.predict(Xte)
112    else:
113        # label transfer by kNN with train labels from full-data 'labels'
114        ytr = np.asarray(labels)[tr]
115        knn = KNeighborsClassifier(n_neighbors=10, weights="distance")
116        knn.fit(Xtr, ytr)
117        yte = knn.predict(Xte)
118
119    return _silhouette_safe(Xte, yte)

```



```

118 # ----- validation core -----
119 def validate_all(
120     X_pca30, X_umap10,
121     K_grid=range(2, 10),
122     dbscan_eps=(0.3, 0.4, 0.5, 0.6, 0.7),
123     dbscan_min_samples=(5, 10, 15),
124     n_boot=10,
125     seed=25
126 ):
127     rng = check_random_state(seed)
128     records = []
129
130     # --- K-based methods ---
131     for K in K_grid:
132         # KMeans (fit/view on UMAP10)
133         km_labels = _kmeans_fit_predict(X_umap10, K, seed)
134         sil = _silhouette_safe(X_umap10, km_labels)
135         stab = _bootstrap_stability(X_umap10, lambda Z: _kmeans_fit_predict(Z,
136                                     K, rng.randint(1e9)), n_boot,
137                                     seed)
138
139         gen = _generalizability_test_silhouette(X_umap10, km_labels, "kmeans",
140                                                 k=K, seed=seed)
141         records.append(("kmeans", K, np.nan, np.nan, sil, stab, gen))
142
143         # GMM (UMAP10)
144         gmm_labels, _ = _gmm_fit_predict(X_umap10, K, seed)
145         sil = _silhouette_safe(X_umap10, gmm_labels)
146         stab = _bootstrap_stability(X_umap10, lambda Z: _gmm_fit_predict(Z, K,
147                                     rng.randint(1e9))[0], n_boot,
148                                     seed)
149
150         gen = _generalizability_test_silhouette(X_umap10, gmm_labels, "gmm", k
151                                                 =K, seed=seed)
152         records.append(("gmm", K, np.nan, np.nan, sil, stab, gen))
153
154         # Spectral (fit on PCA30)
155         sp_labels = _spectral_fit_predict(X_pca30, K, seed)
156         sil = _silhouette_safe(X_pca30, sp_labels)
157         stab = _bootstrap_stability(X_pca30, lambda Z: _spectral_fit_predict(Z
158                                     , K, rng.randint(1e9)), n_boot,
159                                     seed)
160
161         gen = _generalizability_test_silhouette(X_pca30, sp_labels, "spectral"
162                                                 , k=K, seed=seed)
163         records.append(("spectral_pca", K, np.nan, np.nan, sil, stab, gen))
164
165         # Agglomerative (UMAP10): ward/average/complete
166         for lk in ("ward", "average", "complete"):
167             ag_labels = _agglo_fit_predict(X_umap10, K, lk)
168             sil = _silhouette_safe(X_umap10, ag_labels)
169             stab = _bootstrap_stability(X_umap10, lambda Z, _lk=lk:
170                                     _agglo_fit_predict(Z, K,
171                                                         _lk), n_boot, seed)
172
173             gen = _generalizability_test_silhouette(X_umap10, ag_labels, f"
174                                                         agg_{lk}", k=K, seed=seed)
175             records.append((f"agg_{lk}", K, np.nan, np.nan, sil, stab, gen))

```

```

160
161 # --- DBSCAN grid (UMAP10) ---
162 for eps in dbscan_eps:
163     for m in dbscan_min_samples:
164         db_labels = _dbscan_fit_predict(X_umap10, eps, m)
165         sil = _silhouette_safe(X_umap10, db_labels)
166         stab = _bootstrap_stability(
167             X_umap10, lambda Z, _e=eps, _m=m: _dbscan_fit_predict(Z, _e,
168                                                         _m), n_boot, seed
169         )
170         gen = _generalizability_test_silhouette(X_umap10, db_labels, "
171                                                     dbscan", k=None, seed=seed)
172         records.append(("dbscan", np.nan, eps, m, sil, stab, gen))
173
174 cols = ["method", "K", "eps", "min_samples", "silhouette", "stability", "
175         generalizability"]
176 return pd.DataFrame.from_records(records, columns=cols)
177
178 def pick_best_per_method(df: pd.DataFrame) -> pd.DataFrame:
179     out = []
180     for meth in df["method"].unique():
181         sub = df[df["method"] == meth].copy()
182         if sub.empty:
183             continue
184         # rank by (silhouette, stability, generalizability)
185         sub["_rank"] = list(zip(-sub["silhouette"].fillna(-1e9),
186                                 -sub["stability"].fillna(-1e9),
187                                 -sub["generalizability"].fillna(-1e9)))
188         sub = sub.sort_values("_rank").drop(columns=["_rank"])
189         out.append(sub.iloc[[0]])
190     return pd.concat(out, ignore_index=True) if out else pd.DataFrame()
191
192 def run(
193     data_path="Homework 2/Code/Data/BRCA_data.csv",
194     outdir="Homework 2/Code/Result/Problem_2",
195     seed=25
196 ):
197     outdir = Path(outdir); outdir.mkdir(parents=True, exist_ok=True)
198
199     # data + embeddings (reuse 2a)
200     X, _ = load_brca(data_path)
201     # (pca2 not needed here)
202     _, X_pca30, X_umap10, _ = compute_embeddings(X, seed=seed, n_pcs_feat=30)
203
204     # validation
205     df = validate_all(
206         X_pca30, X_umap10,
207         K_grid=range(2, 10),
208         dbscan_eps=(0.3, 0.4, 0.5, 0.6, 0.7),
209         dbscan_min_samples=(5, 10, 15),
210         n_boot=10,
211         seed=seed
212     )

```

```

211 # save all + best-per-method
212 all_csv = outdir / "problem2b_validation.csv"
213 best_csv = outdir / "problem2b_best.csv"
214 df.to_csv(all_csv, index=False)
215 pick_best_per_method(df).to_csv(best_csv, index=False)
216
217 dfk = df[df["K"].notna()].copy()
218 dfk["K"] = dfk["K"].astype(int)
219
220 # unify method labels for legend
221 label_map = {
222     "kmeans": "KMeans",
223     "gmm": "GMM",
224     "spectral_pca": "Spectral (PCA)",
225     "agg_ward": "Agg-Ward",
226     "agg_average": "Agg-Average",
227     "agg_complete": "Agg-Complete",
228 }
229 dfk["Method"] = dfk["method"].map(label_map)
230
231 palette = {
232     "KMeans": "#1f77b4",
233     "GMM": "#2ca02c",
234     "Spectral (PCA)": "#ff7f0e",
235     "Agg-Ward": "#9467bd",
236     "Agg-Average": "#d62728",
237     "Agg-Complete": "#8c564b",
238 }
239
240 # ---- plotting helper ----
241 def plot_metric(metric, ylabel):
242     plt.figure(figsize=(7,5), dpi=130)
243     sns.lineplot(data=dfk, x="K", y=metric, hue="Method", marker="o",
244                 palette=palette)
245     plt.xlabel("Number of Clusters (K)")
246     plt.ylabel(ylabel)
247     plt.title(f"{ylabel} vs K across clustering methods")
248     plt.legend(bbox_to_anchor=(1.02,1), loc="upper left", frameon=False)
249     plt.grid(True, linewidth=0.4, alpha=0.3)
250     plt.tight_layout()
251     plt.savefig(outdir / f"validation_{metric.lower()}_vs_K.png",
252               bbox_inches="tight")
253     plt.close()
254
255 # ---- draw three figures ----
256 plot_metric("silhouette", "Silhouette Score")
257 plot_metric("stability", "Stability (Bootstrap ARI)")
258 plot_metric("generalizability", "Generalizability (Test Silhouette)")
259
260 if __name__ == "__main__":
261     import argparse
262     parser = argparse.ArgumentParser()
263     parser.add_argument("--data_path", type=str, default="Homework 2/Code/Data
264                       /BRCA_data.csv")

```

```

262     parser.add_argument("--outdir", type=str, default="Homework 2/Latex/
263                           Figures")
264     parser.add_argument("--seed", type=int, default=25)
265     args = parser.parse_args()
266     run(data_path=args.data_path, outdir=args.outdir, seed=args.seed)

```

A.5 Problem 2c.

```

1  '''
2  Author: Chuyang Su cs4570@columbia.edu
3  Date: 2025-10-31 02:57:20
4  LastEditTime: 2025-10-31 02:57:22
5  FilePath: /Unsupervised-Learning-Homework/Homework 2/Code/Problem_2_c.py
6  Description:
7      plot the best model.
8  '''
9  import os
10 os.environ.setdefault("OMP_NUM_THREADS", "2")
11 os.environ.setdefault("MKL_NUM_THREADS", "2")
12
13 import warnings
14 warnings.filterwarnings("ignore")
15
16 from pathlib import Path
17 import numpy as np
18 import pandas as pd
19 import matplotlib.pyplot as plt
20 import seaborn as sns
21
22 from sklearn.cluster import AgglomerativeClustering
23 from sklearn.metrics import adjusted_rand_score, normalized_mutual_info_score
24
25 from Problem_2_a import load_brca, compute_embeddings, scatter_by_labels
26
27 def best_clustering_labels(X_umap10, k=2):
28     model = AgglomerativeClustering(n_clusters=k, linkage="average")
29     return model.fit_predict(X_umap10)
30
31 def run(
32     data_path="Homework 2/Code/Data/BRCA_data.csv",
33     outdir="Homework 2/Code/Result/Problem_2",
34     seed=25
35 ):
36     outdir = Path(outdir)
37     outdir.mkdir(parents=True, exist_ok=True)
38
39     X, meta = load_brca(data_path)
40     X_pca2, X_pca30, X_umap10, X_umap2 = compute_embeddings(X, seed=seed,
41                                                             n_pcs_feat=30)
42
43     pam50 = meta["subtype"].astype(str).fillna("NA").values
44     labels_best = best_clustering_labels(X_umap10, k=2)

```

```

44
45 ari = adjusted_rand_score(pam50, labels_best)
46 nmi = normalized_mutual_info_score(pam50, labels_best)
47 print(f"ARI = {ari:.3f}, NMI = {nmi:.3f}")
48
49 # Define consistent color map
50 PAM50_COLOR_MAP = {
51     "Luminal A": "#1f77b4",
52     "Luminal B": "#2ca02c",
53     "Basal-like": "#ff7f0e",
54     "HER2-enriched": "#d62728",
55     "Normal-like": "#9467bd",
56 }
57
58 # Best clustering (Agg-Average, K=2)
59 cmap_best = {str(i): c for i, c in enumerate(["#4C78A8", "#E45756"])}
60 scatter_by_labels(
61     X_umap2, labels_best.astype(str),
62     title="Best Clustering (Agglomerative Average, K=2)",
63     outpath=outdir / "2c_bestcluster_umap2.png",
64     color_map=cmap_best,
65     legend_title="Cluster"
66 )
67
68 # Combined overlay (color by PAM50, edge by cluster)
69 df_plot = pd.DataFrame({
70     "x": X_umap2[:,0], "y": X_umap2[:,1],
71     "Cluster": labels_best.astype(str),
72     "Subtype": pam50
73 })
74
75 plt.figure(figsize=(7,5.5), dpi=130)
76 sns.scatterplot(
77     data=df_plot, x="x", y="y",
78     hue="Subtype", style="Cluster",
79     palette=PAM50_COLOR_MAP, alpha=0.9, s=28
80 )
81 plt.title(f"Overlay of PAM50 and Best Clustering (ARI={ari:.2f}, NMI={nmi:.2f})")
82 plt.xlabel("UMAP-1"); plt.ylabel("UMAP-2")
83 plt.legend(bbox_to_anchor=(1.02,1), loc="upper left", frameon=False)
84 plt.grid(True, linewidth=0.3, alpha=0.3)
85 plt.tight_layout()
86 plt.savefig(outdir / "2c_overlay_umap2.png", bbox_inches="tight")
87 plt.close()
88
89 if __name__ == "__main__":
90     import argparse
91     parser = argparse.ArgumentParser()
92     parser.add_argument("--data_path", type=str, default="Homework 2/Code/Data/BRCA_data.csv")
93     parser.add_argument("--outdir", type=str, default="Homework 2/Latex/Figures")
94     parser.add_argument("--seed", type=int, default=25)

```

```
95     args = parser.parse_args()
96     run(data_path=args.data_path, outdir=args.outdir, seed=args.seed)
```