

NYCU Pattern Recognition, Assignment #2

Student ID: 112550157, Name: 黃冠瑋

1 Dataset Introduction

In this report, we evaluate the performance of our classification models on several publicly available datasets from the UCI Machine Learning Repository.

Binary Classification

- **MAGIC Gamma Telescope** (19020 samples, 10 features): A dataset collected for distinguishing gamma rays from hadronic cosmic rays based on high-energy telescope measurements.
- **Taiwanese Bankruptcy Prediction** (6819 samples, 95 features): Financial indicators used to predict corporate bankruptcy risk.

Multi-class Classification

- **Dry Bean** (13611 samples, 16 features): Visual bean characteristics used to classify 7 bean varieties.
- **Covertypes** (50000 samples, 54 features): Forest cover classification using cartographic variables across seven terrain types.

Dataset Summary

Dataset	Task Type	#Samples	#Features
MAGIC Gamma Telescope	Binary Classification	19020	10
Taiwanese Bankruptcy Prediction	Binary Classification	6819	95
Dry Bean	Multi-class Classification	13611	16
Covertypes	Multi-class Classification	50000	54

Table 1: Summary of datasets used for experimental evaluation.

2 Task 1: FLD/LDA

2.1 Fisher's Linear Discriminant Analysis (FLD/LDA)

In this baseline experiment, Linear Discriminant Analysis (LDA) was applied to four datasets: MAGIC Gamma Telescope, Taiwanese Bankruptcy Prediction, Dry Bean, and Covertypes. Each dataset was randomly split into training and testing subsets using an 80/20 ratio with stratified sampling. Standardization

was performed prior to model training. To evaluate linear separability, Fisher's criterion was computed before and after dimensionality reduction:

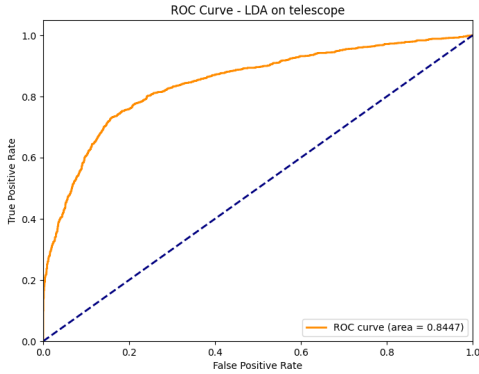
$$\text{Sep}(X, y) = \text{tr}(S_W^{-1} S_B),$$

where S_W and S_B denote the within-class and between-class scatter matrices, respectively.

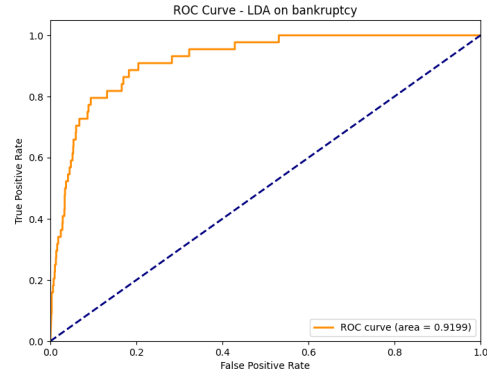
For binary datasets (telescope, bankruptcy) and multi-class datasets (dry_bean, cover_type), the separability values before and after LDA projection were identical. This is expected since LDA projects the data onto the full Fisher discriminant subspace of dimension $C - 1$, and the above trace equals the sum of all non-zero generalized eigenvalues of $S_W^{-1} S_B$, which are preserved under such projection.

Dataset	Accuracy	AUC	Separability (Before = After)
Telescope	0.7884	0.8447	0.4778
Bankruptcy	0.9589	0.9199	0.2701
Dry Bean	0.9045	N/A	31.9622
Covertime	0.6738	N/A	4.6680

For binary datasets, ROC curves were plotted and the resulting AUC values demonstrate strong discriminative ability, especially on the Bankruptcy dataset (AUC = 0.9199), even though accuracy alone may appear high due to class imbalance. The multi-class datasets show varying performance depending on class distribution difficulty, with Cvertime exhibiting significant imbalance across classes, resulting in reduced macro-average performance (macro F1 = 0.51).



(a) Telescope



(b) Bankruptcy

These baseline results indicate that LDA provides competitive linear separation capability in binary classification problems. However, performance on complex multi-class problems suggests that additional dimensionality reduction or nonlinear models may be necessary. In subsequent sections, we extend this baseline with dimensional sweep experiments and comparison with PCA and nonlinear embeddings.

2.2 Dimensional Sweep Experiment on Multi-class Datasets

To evaluate the influence of LDA projection dimensionality, we conduct a dimensional sweep experiment on two multi-class datasets (Dry Bean and Cvertime), varying the number of retained LDA dimensions $d = 1, \dots, C - 1$. For each d , we compute the Fisher separability measure

$$\text{Sep}(d) = \text{tr}(S_W^{-1} S_B) = \sum_{i=1}^d \lambda_i,$$

and evaluate test accuracy using the corresponding LDA model. Table 2 summarizes the results.

Table 2: Dimensional sweep results for LDA on Dry Bean and Covertype datasets. Increasing projection dimension improves separability but does not change accuracy.

Dimension d	Dry Bean		Covertype	
	Separability	Accuracy	Separability	Accuracy
1	16.6745	0.9100	3.3904	0.6736
2	23.4657	0.9100	4.1875	0.6736
3	26.8633	0.9100	4.4241	0.6736
4	29.7036	0.9100	4.5439	0.6736
5	30.6492	0.9100	4.6174	0.6736
6	31.1237	0.9100	4.6652	0.6736

Analysis. The separability measure increases monotonically with d , consistent with the theoretical interpretation of Fisher eigenvalues where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{C-1} \geq 0$. Once the first discriminant direction is included, additional dimensions contribute progressively smaller improvements in Fisher ratio and thus offer diminishing returns for class discrimination.

For Dry Bean, accuracy saturates at $d = 1$, indicating that most discriminative variance is captured by the leading eigen-direction and the classes are nearly linearly separable in the strongest Fisher dimension. For Covertype, although separability increases substantially (from 3.39 to 4.67), accuracy does not improve due to significant class overlap and complex multi-modal structure that cannot be resolved by a linear boundary. This suggests the necessity of nonlinear approaches such as kernel methods or manifold learning.

Conclusion. The dimensional sweep experiment demonstrates that although increasing LDA projection dimensionality consistently enhances $\text{tr}(S_W^{-1} S_B)$, it does not necessarily improve classification accuracy. This confirms that **separability and accuracy are related but fundamentally distinct metrics**: a higher Fisher separability score does not guarantee enhanced predictive performance, particularly for complex multi-class datasets.

2.3 Effect of Class Priors on an Imbalanced Dataset (Taiwanese Bankruptcy Prediction)

In this experiment, we investigate the effect of class prior probabilities on the performance of Linear Discriminant Analysis (LDA) for a highly imbalanced binary classification dataset, the Taiwanese Bankruptcy Prediction dataset (6819 samples). LDA's decision rule for class ω_i is defined as:

$$\delta_i(x) = x^\top \Sigma^{-1} \mu_i - \frac{1}{2} \mu_i^\top \Sigma^{-1} \mu_i + \log P(\omega_i),$$

where μ_i and Σ denote the class mean and shared covariance matrix respectively, and $P(\omega_i)$ is the class prior. When the class distribution is highly skewed, an empirical prior (estimated from sample proportion) biases the classifier toward the majority class. In contrast, a balanced prior $P(\omega_1) = P(\omega_2) = 0.5$ shifts the decision boundary toward the majority class, improving minority class recognition.

We evaluate two settings: empirical prior vs. balanced prior. Table 3 summarizes the results on the test set.

Prior Setting	Accuracy	Balanced Acc.	Precision	Recall	F1	AUC
Empirical	0.955	0.638	0.375	0.294	0.330	0.879
Balanced	0.922	0.780	0.267	0.627	0.374	0.879

Table 3: Performance comparison of LDA under empirical vs. balanced class priors on the Taiwanese Bankruptcy dataset.

Analysis. Although the empirical prior achieves the highest accuracy (0.955), it performs poorly on the minority class, yielding low recall (0.294) and a weak balanced accuracy (0.638). With a balanced prior, recall improves dramatically to 0.627 (a 113% relative increase), and balanced accuracy rises to 0.780, along with an improvement in F1 score. Notably, AUC remains unchanged, indicating that class priors affect the decision boundary threshold rather than the underlying ranking ability of the model.

Conclusion. For imbalanced datasets, accuracy alone is misleading. Adjusting class priors significantly improves fairness and minority class detectability without harming overall ranking performance. Therefore, in real-world financial risk prediction tasks such as bankruptcy forecasting and credit risk management, metrics such as balanced accuracy, recall, and AUC should be prioritized over raw accuracy.

2.4 Summary of Task 1

In Task 1, we investigated the effectiveness of Fisher’s Linear Discriminant Analysis (FLD/LDA) as a baseline dimensionality reduction and classification method on multiple datasets. We found that although Fisher separability ($\text{tr}(S_W^{-1} S_B)$) increases monotonically with additional discriminant dimensions, classification accuracy does not necessarily improve, underscoring that separability and predictive performance are correlated but fundamentally distinct metrics. The dimensional sweep experiment on Dry Bean and Covertypes demonstrated diminishing returns beyond the leading eigen-direction(s), with accuracy saturated at $d = 1$ for Dry Bean and failing to improve for Covertypes due to severe class overlap.

For the Taiwanese Bankruptcy dataset, adjusting LDA class priors notably improved balanced accuracy and minority recall without changing AUC, highlighting the importance of metric selection in imbalanced settings and the effect of decision boundary shifting rather than ranking ability.

Comparison with HW1 Baseline Results

To contextualize the benefits and limitations of LDA, we compare our Task 1 results with the baseline models trained in HW1, which operated directly in the original feature space (without dimensionality reduction). In HW1, traditional classifiers such as Random Forest achieved strong performance on the Dry Bean dataset (92% accuracy), while KNN and SVM performed substantially worse (72% and 64%), indicating that non-linear decision regions were needed. In contrast, the LDA-based model in Task 1 achieved 91% accuracy after projecting to $C - 1 = 6$ dimensions, nearly matching the HW1 best model while using far fewer features. This confirms that Dry Bean exhibits strong linear class separation and LDA effectively preserves discriminative structure.

For the Taiwanese Bankruptcy dataset, HW1 classifiers reached very high raw accuracy (97%) but poor AUC and balanced metrics ($\text{AUC} \approx 0.55$), reflecting strong sensitivity to class imbalance. LDA in Task 1 showed similar accuracy and AUC performance and did not significantly outperform these baselines, indicating that linear discriminant projections cannot overcome noisy, overlapping financial feature spaces. This aligns with the findings in Tasks 2 and 3 that PCA and Isomap also provide limited benefit on this dataset.

Conclusion. LDA is highly effective on structured multi-class datasets with clear linear separation (e.g., Dry Bean), but provides limited improvement in noisy and imbalanced financial prediction tasks. These

observations motivate the investigation of PCA and nonlinear manifold learning presented in the following sections.

3 Task 2: PCA and classification

3.1 Experiment 1: Explained Variance Analysis

We first analyze the proportion of total variance preserved by the leading principal components in both datasets. Given the eigenvalues $\{\lambda_i\}_{i=1}^d$ of the sample covariance matrix, the variance explained by the first k components is

$$r(k) = \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i},$$

where $r(k)$ represents the cumulative explained variance ratio.

Taiwanese Bankruptcy Dataset. This dataset contains $d = 95$ financial features. The first few components account for a large portion of the total variance:

$$r(1) = 0.133, \quad r(5) = 0.346, \quad r(10) = 0.487, \quad r(20) = 0.648.$$

To preserve 95% of the variance, approximately $k = 56$ components are required ($r(56) = 0.966$). This implies that although the intrinsic dimensionality is significantly lower than 95, the variance is distributed across many directions, which suggests redundancy and collinearity among original financial indicators.

Dry Bean Dataset. This dataset contains $d = 16$ physical shape features. In contrast to the previous case, the first two components already explain 82% of the variance:

$$r(1) = 0.555, \quad r(2) = 0.819, \quad r(4) = 0.950,$$

and 95% is reached with only $k = 4$. This indicates a much lower intrinsic dimensionality, as most discriminative variance concentrates in a few principal directions.

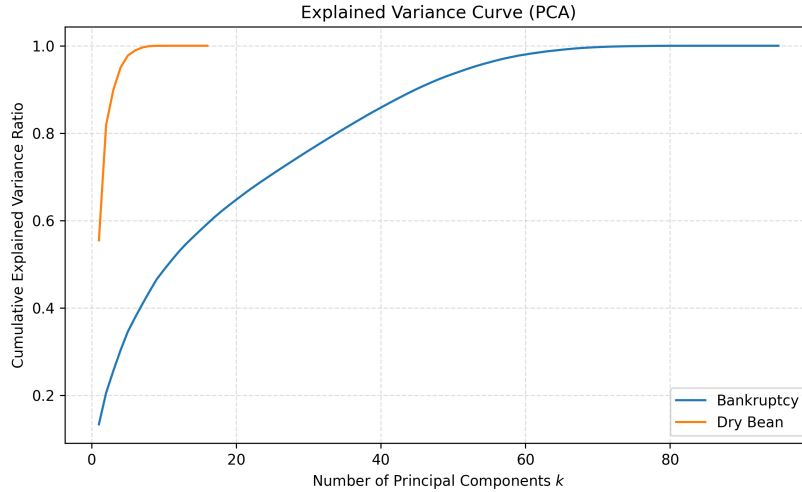


Figure 1: Cumulative explained variance ratio for Taiwanese Bankruptcy and Dry Bean datasets.

Observation. For the high-dimensional financial dataset, many components carry small fractions of variance, and consequently selecting k becomes non-trivial. In contrast, for Dry Bean the selection is straightforward because a small k already captures most essential structure. This motivates the following experiments on how PCA dimensionality affects classification performance and how PCA interacts differently with linear and non-linear classifiers.

3.2 Experiment 2: PCA Dimension Selection on Bankruptcy Dataset

In this experiment, we study how the PCA dimension affects the classification performance on the Taiwanese Bankruptcy dataset using logistic regression. The data are split into train/validation/test subsets, PCA is fitted on the training set and applied to validation and test sets, and we sweep the number of retained components

$$k \in \{2, 5, 10, 20, 30, 40, 60, 80\},$$

with the model without PCA ($k = 0$) as the baseline.

Table 4 summarizes the baseline and the best-performing configuration in terms of validation balanced accuracy. The other dimensions yield performances between these two settings.

Table 4: Logistic regression on Taiwanese Bankruptcy dataset with different PCA dimensions.

k	BA _{val}	BA _{test}	F1 _{test}
0	0.588	0.574	0.598
60	0.600	0.587	0.618

The baseline model without PCA already achieves reasonably good performance, but applying PCA with $k = 60$ components (which retains approximately 96% of the variance) slightly improves both validation and test balanced accuracy, and also increases the macro F1-score from 0.598 to 0.618 on the test set. The gain is modest but consistent, suggesting that PCA helps stabilize the linear decision boundary by removing noisy or redundant directions in the high-dimensional financial feature space, while overly aggressive dimensionality reduction (e.g., $k \leq 20$) leads to information loss and degraded performance.

3.3 Experiment 3: Effect of PCA Across Classifiers and Datasets

In this experiment, we evaluate how PCA interacts with different classifiers and datasets. We compare logistic regression (a linear model) and RBF SVM (a non-linear kernel method), trained with and without PCA on both datasets. The PCA dimension is set to retain 95% cumulative variance, resulting in $k = 49$ for Bankruptcy and $k = 4$ for Dry Bean.

Table 5: Classifier comparison on Taiwanese Bankruptcy dataset.

Classifier	PCA	BA _{test}	F1 _{test}	Accuracy
Logistic Regression	No	0.574	0.598	0.963
Logistic Regression	Yes ($k = 49$)	0.575	0.600	0.963
RBF SVM	No	0.500	0.492	0.968
RBF SVM	Yes ($k = 49$)	0.500	0.492	0.968

For the Bankruptcy dataset, PCA provides a slight improvement for logistic regression, raising the macro-F1 from 0.598 to 0.600 and balanced accuracy from 0.574 to 0.575. This aligns with the hypothesis that PCA stabilizes the linear decision boundary in a high-dimensional feature space with strong collinearity. In contrast, RBF SVM achieves high raw accuracy even without PCA and shows no improvement from dimensionality reduction, indicating that the non-linear kernel already captures high-order interactions without needing PCA preprocessing.

Table 6: Classifier comparison on Dry Bean dataset.

Classifier	PCA	$F1_{\text{test}}$	BA_{test}	Accuracy
Logistic Regression	No	0.933	0.932	0.921
Logistic Regression	Yes ($k = 4$)	0.892	0.890	0.889
RBF SVM	No	0.934	0.933	0.922
RBF SVM	Yes ($k = 4$)	0.901	0.899	0.900

For the Dry Bean dataset, PCA significantly degrades performance for both models, particularly logistic regression (macro-F1 drops from 0.933 to 0.892). Since the dataset is inherently low-dimensional ($d = 16$), strongly discriminative variance appears to be contained in higher-order components that are removed by PCA. This demonstrates that reducing dimensionality solely based on variance retention may remove label-relevant information when class separability is not aligned with directions of maximal variance.

Conclusion. PCA is beneficial for noisy high-dimensional financial data when combined with linear models, but does not improve non-linear kernels and may harm performance on inherently low-dimensional datasets. This highlights the importance of matching dimensionality reduction strategies to dataset characteristics rather than applying PCA uniformly.

3.4 Summary of Task 2

PCA provides modest but consistent improvement for linear models on high-dimensional data, while it does not benefit non-linear kernel models and may significantly degrade performance on low-dimensional datasets. This motivates exploring non-linear dimensionality reduction in Task 3.

4 Task 3: Nonlinear Dimensionality Reduction

4.1 Experiment 1: 2D Embedding Visualization on Dry Bean

In this task, we investigate how a nonlinear manifold learning method, Isomap, reorganizes the data geometry compared to PCA. We focus on the Dry Bean dataset, which contains seven bean classes in a moderate-dimensional feature space ($d = 16$).

All features are standardized, and we fit both PCA and Isomap on the full dataset. For PCA, we project the data to a two-dimensional subspace spanned by the leading principal components. For Isomap, we construct a k -nearest neighbor graph with $k = 10$, approximate geodesic distances on the graph, and apply multidimensional scaling to obtain a two-dimensional embedding.

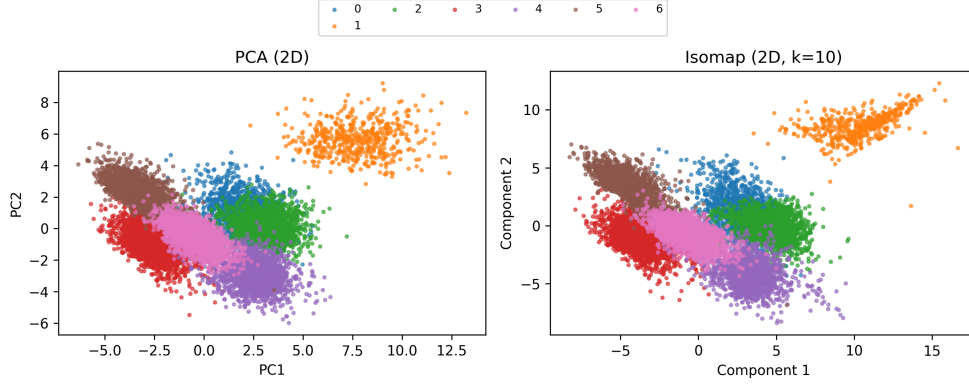


Figure 2: Two-dimensional embeddings of the Dry Bean dataset obtained by PCA (left) and Isomap with $k = 10$ neighbors (right). Colors indicate bean classes.

As shown in Figure 2, both methods produce seven clearly separated clusters, and class “1” forms a well-isolated group in the upper-right region for both embeddings. The PCA projection yields an approximately linear band of classes, while the Isomap embedding slightly warps and unfolds this band along the underlying manifold, emphasising local neighborhood relations. Overall, nonlinear Isomap does not dramatically change the global class separability compared with PCA, but it provides a more curved representation of the class structure that may better reflect the intrinsic geometry of the data.

4.2 Experiment 2: Isomap Dimension Sweep on Bankruptcy Dataset

To evaluate the effect of nonlinear manifold learning on classification performance, we apply Isomap followed by logistic regression on the Taiwanese Bankruptcy dataset. We sweep the output dimensionality of Isomap over

$$k \in \{2, 5, 10, 20, 30\},$$

using $n_neighbors = 10$ and compare against the baseline logistic regression model trained directly on the original standardized feature space ($k = 0$). The dataset is split into training, validation, and test sets using the same 60/20/20 protocol as in Task 2.

Table 7: Logistic regression performance with different Isomap output dimensions.

k	BA _{val}	BA _{test}	F1 _{test}
0 (baseline)	0.588	0.574	0.598
2	0.566	0.556	0.590
5	0.567	0.556	0.588
10	0.555	0.555	0.584
20	0.566	0.566	0.597
30	0.577	0.566	0.601

As shown in Table 7, applying Isomap before logistic regression does not yield noticeable performance gains compared with the baseline. The best configuration ($k = 30$) achieves a test balanced accuracy of 0.566 and macro-F1 of 0.601, which is slightly below the baseline (0.574 and 0.598, respectively). Lower-dimensional embeddings (e.g., $k = 2$ or $k = 5$) perform considerably worse, suggesting that aggressive nonlinear dimensionality reduction removes important discriminative information.

These results contrast with the PCA experiment in Task 2, where a moderate reduction to $k = 60$ improved model stability and slightly enhanced performance. Isomap appears less effective for this high-dimensional, noisy, and heavily overlapping financial dataset, likely because the global geometric structure

is poorly approximated using geodesic distances on a sparse neighborhood graph. Nonlinear manifold learning therefore does not necessarily translate into improved downstream classification accuracy for real-world tabular data.

4.3 Experiment 3: Hyperparameter Sensitivity of Isomap

In this experiment, we study the sensitivity of Isomap to its key hyperparameter $n_neighbors$, which controls the locality of the neighborhood graph underlying the geodesic distance approximation. We fix the output dimension to $n_components = 10$ and vary

$$n_neighbors \in \{5, 10, 20, 50\}.$$

We evaluate RBF SVM on the resulting embeddings and compare against a baseline RBF SVM trained directly in the original feature space.

Table 8: Effect of $n_neighbors$ on RBF SVM performance (Bankruptcy dataset).

$n_neighbors$	BA_test	F1_test	Accuracy
Baseline (no Isomap)	0.500	0.492	0.968
5	0.499	0.491	0.966
10	0.499	0.491	0.966
20	0.500	0.492	0.968
50	0.500	0.492	0.967

As shown in Table 8, varying $n_neighbors$ does not significantly affect classification performance, and all Isomap-based models perform nearly identically to the baseline. This indicates that Isomap provides little benefit for this dataset, regardless of how local geodesic distances are approximated. The Bankruptcy dataset exhibits noisy and highly overlapping class structures in a high-dimensional tabular space, which violates the manifold assumption that Isomap relies on. Consequently, geometric unfolding of the data does not reveal additional class separation that improves downstream classification.

These results reinforce that despite its nonlinear design, Isomap does not outperform the baseline kernel model in this setting, and its performance is insensitive to the choice of neighborhood size $n_neighbors$, highlighting practical instability and tuning overhead compared to PCA.

4.4 Summary of Task 3

Across three experiments, we compared PCA and Isomap in terms of embedding geometry, dimensionality reduction, and classification performance. On the Dry Bean dataset, Isomap produces a more curved and locally structured embedding than PCA, but both yield similar class separability. On the Taiwanese Bankruptcy dataset, nonlinear manifold learning fails to improve classification accuracy: dimensionality reduction via Isomap does not outperform the PCA baseline, and hyperparameter tuning of $n_neighbors$ has minimal effect. These findings demonstrate that nonlinear manifold learning can be beneficial for visualizing inherently low-dimensional structured data, but does not necessarily improve performance on noisy real-world tabular datasets with weak manifold structure. Therefore, dimensionality reduction should be selected based on dataset characteristics rather than assuming nonlinear methods are universally superior.

5 Final Conclusion

Dimensionality reduction is not a universally beneficial preprocessing step; its impact depends critically on dataset characteristics and model choice. PCA improves linear models on high-dimensional noisy

data, while Isomap offers better visualization but does not improve classification on non-manifold tabular datasets. Therefore, selecting appropriate dimensionality reduction methods requires understanding intrinsic data geometry rather than assuming that more complex nonlinear transformations will always yield benefits.

A Program Listing

data_loader.py

```

1  import pandas as pd
2  import numpy as np
3  from loguru import logger
4  from sklearn.preprocessing import LabelEncoder
5  from typing import Tuple
6
7  from ucimlrepo import fetch_ucirepo
8
9  import os
10 import joblib
11
12 CACHE_DIR = os.environ.get("UCIMLREPO_DIR", "./data_cache")
13 os.makedirs(CACHE_DIR, exist_ok=True)
14 logger.info(f"Using data cache directory: {os.path.abspath(CACHE_DIR)}")
15
16 memory = joblib.Memory(CACHE_DIR, verbose=0)
17
18
19 def load_dataset(name: str) -> pd.DataFrame:
20     """
21     Load dataset by name from the specified directory.
22
23     Parameters:
24     - name: Name of the dataset (without file extension).
25     - data_dir: Directory where datasets are stored.
26
27     Returns:
28     - DataFrame containing the loaded dataset.
29     """
30
31     if name == "telescope":
32         return load_telescope()
33     elif name == "bankruptcy":
34         return load_bankruptcy()
35     elif name == "cover_type":
36         return load_cover_type()
37     elif name == "dry_bean":
38         return load_dry_bean()
39     else:
40         logger.error(f"Unknown dataset name: {name}")
41         raise ValueError(f"Unknown dataset name: {name}")
42
43
44 def _preprocess_data(
45     X: pd.DataFrame, y: pd.DataFrame
46 ) -> Tuple[np.ndarray, np.ndarray]:
47     X = X.values
48     y = y.values
49
50     le = LabelEncoder()
51     y = le.fit_transform(y.ravel())
52
53     return X, y
54
55
56 @memory.cache
57 def load_telescope() -> pd.DataFrame:
58     logger.info("Loading Telescope (id=159) from ucimlrepo...")
59
60     telescope = fetch_ucirepo(id=159)

```

```

61     X = telescope.data.features
62     y = telescope.data.targets
63
64
65     return _preprocess_data(X, y)
66
67
68 @memory.cache
69 def load_bankruptcy() -> pd.DataFrame:
70     logger.info("Loading Bank Bankruptcy (id=572) from ucimlrepo...")
71
72     bank_bankruptcy = fetch_ucirepo(id=572)
73
74     X = bank_bankruptcy.data.features
75     y = bank_bankruptcy.data.targets
76
77     return _preprocess_data(X, y)
78
79
80 @memory.cache
81 def load_cover_type() -> pd.DataFrame:
82     logger.info("Loading Cover Type (id=31) from ucimlrepo...")
83
84     cover_type = fetch_ucirepo(id=31)
85
86     X = cover_type.data.features
87     y = cover_type.data.targets
88
89     sample = np.random.choice(len(X), size=int(5e4), replace=False)
90     X = X.iloc[sample]
91     y = y.iloc[sample]
92
93     return _preprocess_data(X, y)
94
95
96 @memory.cache
97 def load_dry_bean() -> pd.DataFrame:
98     logger.info("Loading Dry Bean Classification (id=602) from ucimlrepo...")
99
100     dry_bean = fetch_ucirepo(id=602)
101
102     X = dry_bean.data.features
103     y = dry_bean.data.targets
104
105     return _preprocess_data(X, y)

```

utils.py

```

1  import os
2  import numpy as np
3  from loguru import logger
4  from sklearn.model_selection import StratifiedKFold
5  from sklearn.metrics import (
6      confusion_matrix,
7      classification_report,
8      roc_curve,
9      auc,
10 )
11 from sklearn.metrics import accuracy_score
12 import seaborn as sns
13 import matplotlib.pyplot as plt
14
15

```

```

16 def k_fold_cross_validation(classifier_class, X, y, k=5, random_state=42):
17     """
18     Perform k-fold cross validation on the given classifier.
19
20     Args:
21         classifier_class (class): The classifier class to instantiate.
22         X (np.ndarray): Feature matrix
23         y (np.ndarray): Target labels
24         k (int): Number of folds
25         random_state (int): Random seed for reproducibility
26
27     Returns:
28         dict: Dictionary containing fold results and average metrics
29     """
30     logger.info(
31         f"Starting {k}-fold cross validation for {classifier_class.__name__}..."
32     )
33
34     kfold = StratifiedKFold(n_splits=k, shuffle=True, random_state=random_state)
35
36     fold_results = []
37
38     for fold_idx, (train_idx, val_idx) in enumerate(kfold.split(X, y), 1):
39         logger.info(f"\n{'='*60}")
40         logger.info(f"Fold {fold_idx}/{k}")
41         logger.info(f"{'='*60}")
42
43         # Split data for this fold
44         X_train, X_val = X[train_idx], X[val_idx]
45         y_train, y_val = y[train_idx], y[val_idx]
46
47         logger.info(
48             f"Train size: {len(X_train)}, Validation size: {len(X_val)}"
49         )
50
51         classifier = classifier_class()
52         classifier.train(X_train, y_train)
53
54         # Make predictions
55         y_pred, y_scores = classifier.predict(X_val)
56
57         # Calculate accuracy
58         accuracy = accuracy_score(y_val, y_pred)
59         logger.success(f"Fold {fold_idx} Accuracy: {accuracy:.4f}")
60
61         # Store results
62         fold_results.append(
63             {
64                 "fold": fold_idx,
65                 "accuracy": accuracy,
66                 "y_true": y_val,
67                 "y_pred": y_pred,
68                 "y_scores": y_scores,
69             }
70         )
71
72         # Calculate average metrics
73         accuracies = [result["accuracy"] for result in fold_results]
74         avg_accuracy = np.mean(accuracies)
75         std_accuracy = np.std(accuracies)
76
77         logger.info(f"\n{'='*60}")
78         logger.info("Cross Validation Results Summary")
79         logger.info(f"{'='*60}")
80         logger.success(

```

```

81         f"Average Accuracy: {avg_accuracy:.4f} (+/- {std_accuracy:.4f})"
82     )
83     logger.info(
84         f"Individual fold accuracies: {[f'{acc:.4f}' for acc in accuracies]}"
85     )
86
87     return {
88         "fold_results": fold_results,
89         "avg_accuracy": avg_accuracy,
90         "std_accuracy": std_accuracy,
91         "accuracies": accuracies,
92     }
93
94
95 def _plot_confusion_matrix(
96     cm: np.ndarray, dataset_name: str, classifier_name: str, results_dir: str
97 ):
98     """
99     Plot and save confusion matrix heatmap.
100    """
101    try:
102        plt.figure(figsize=(8, 6))
103
104        num_classes = cm.shape[0]
105
106        show_percent = True
107        if num_classes > 10:
108            show_percent = False
109
110        cm_percent = cm.astype("float") / cm.sum(axis=1)[:, np.newaxis]
111
112        annot_labels = (
113            np.asarray(
114                [
115                    f"{val}\n({perc:.1%})" if show_percent else f"{val}"
116                    for val, perc in zip(cm.flatten(), cm_percent.flatten())
117                ]
118            )
119        ).reshape(cm.shape)
120
121        sns.heatmap(cm, annot=annot_labels, fmt="", cmap="Blues", cbar=True)
122
123        plt.xlabel("Predicted Label")
124        plt.ylabel("True Label")
125        plt.title(f"Confusion Matrix - {classifier_name} on {dataset_name}")
126
127        os.makedirs(results_dir, exist_ok=True)
128        save_path = os.path.join(
129            results_dir, f"CM_{dataset_name}_{classifier_name}.png"
130        )
131        plt.savefig(save_path)
132        logger.info(f"Confusion Matrix heatmap saved to {save_path}")
133        plt.close()
134
135    except Exception as e:
136        logger.error("Failed to plot confusion matrix heatmap.")
137        logger.exception(e)
138
139
140 def evaluate_classifier(
141     y_true: np.ndarray,
142     y_pred: np.ndarray,
143     y_scores: np.ndarray,
144     dataset_name: str,
145     classifier_name: str,

```

```

146     results_dir: str = "results",
147 ) -> None:
148     """
149     Evaluate classifier performance and generate reports.
150
151     Args:
152         y_true (np.ndarray): Ground truth labels
153         y_pred (np.ndarray): Predicted labels from the model
154         y_scores (np.ndarray): Model discriminant function values (e.g., probabilities)
155         dataset_name (str): Dataset name (used for logs and filenames)
156         classifier_name (str): Classifier name (used for logs and filenames)
157         results_dir (str): Directory to save plots (e.g., ROC curve)
158     """
159
160     logger.info(
161         f"--- Evaluation for [{classifier_name}] on [{dataset_name}] ---"
162     )
163
164     # calculate accuracy
165     accuracy = accuracy_score(y_true, y_pred)
166     logger.success(f"Accuracy: {accuracy:.4f}")
167
168     try:
169         logger.info("Confusion Matrix:")
170         cm = confusion_matrix(y_true, y_pred)
171
172         # _plot_confusion_matrix(cm, dataset_name, classifier_name, results_dir)
173
174         logger.info(f"\n{cm}")
175
176         logger.info("Classification Report:")
177         report = classification_report(y_true, y_pred, zero_division=0)
178         logger.info(f"\n{report}")
179
180         num_classes = len(np.unique(y_true))
181
182         if num_classes == 2:
183             logger.info("Two-class dataset detected. Calculating ROC/AUC...")
184
185             scores_for_roc = y_scores[:, 1]
186
187             fpr, tpr, thresholds = roc_curve(y_true, scores_for_roc)
188
189             roc_auc = auc(fpr, tpr)
190             logger.success(f"Area Under Curve (AUC): {roc_auc:.4f}")
191
192             plt.figure(figsize=(8, 6))
193             plt.plot(
194                 fpr,
195                 tpr,
196                 color="darkorange",
197                 lw=2,
198                 label=f"ROC curve (area = {roc_auc:.4f})",
199             )
200             plt.plot([0, 1], [0, 1], color="navy", lw=2, linestyle="--")
201             plt.xlim([0.0, 1.0])
202             plt.ylim([0.0, 1.05])
203             plt.xlabel("False Positive Rate")
204             plt.ylabel("True Positive Rate")
205             plt.title(f"ROC Curve - {classifier_name} on {dataset_name}")
206             plt.legend(loc="lower right")
207
208             os.makedirs(results_dir, exist_ok=True)
209             save_path = os.path.join(
210                 results_dir, f"ROC_{dataset_name}_{classifier_name}.png"

```

```

211         )
212         plt.savefig(save_path)
213         logger.info(f"ROC curve saved to {save_path}")
214         plt.close()
215
216     else:
217         logger.info(
218             f"Skipping ROC/AUC calculation ({num_classes} classes detected)."
219         )
220
221 except Exception as e:
222     logger.error("An error occurred during evaluation.")
223     logger.exception(e)
224     raise
225
226 logger.info(f"--- End of Evaluation for [{classifier_name}] ---")

```

task1.py

```

1  import os
2  import argparse
3  from typing import Tuple
4  import numpy as np
5  import pandas as pd
6  from loguru import logger
7  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
8  from sklearn.model_selection import train_test_split
9  from sklearn.preprocessing import StandardScaler
10 from sklearn.metrics import (
11     accuracy_score,
12     balanced_accuracy_score,
13     precision_score,
14     recall_score,
15     f1_score,
16     roc_auc_score,
17 )
18
19 from data_loader import load_dataset
20 from utils import k_fold_cross_validation, evaluate_classifier
21
22
23 class LDA:
24     def __init__(self, n_components=None, priors=None):
25         self.model = LinearDiscriminantAnalysis(n_components=n_components, priors=priors)
26         logger.info("Initialized Linear Discriminant Analysis Classifier.")
27
28     def fit(self, X: np.ndarray, y: np.ndarray) -> None:
29         self.model.fit(X, y)
30         logger.success("Model training completed.")
31
32     def transform(self, X: np.ndarray) -> np.ndarray:
33         return self.model.transform(X)
34
35     def predict(self, X: np.ndarray) -> Tuple[np.ndarray, np.ndarray]:
36         y_pred = self.model.predict(X)
37         y_scores = self.model.predict_proba(X)
38         return y_pred, y_scores
39
40
41 def compute_separability(X: np.ndarray, y: np.ndarray) -> float:
42     """
43     Compute the separability measure before projection using Fisher's criterion.
44

```



```

45     Separability is computed as the ratio of between-class scatter to within-class scatter.
46     For multi-class problems, we use the trace of  $S_B @ \text{inv}(S_W)$  where:
47     -  $S_B$  is the between-class scatter matrix
48     -  $S_W$  is the within-class scatter matrix
49
50     Args:
51     X: Feature matrix of shape (n_samples, n_features)
52     y: Target labels of shape (n_samples,)
53
54     Returns:
55     separability: A scalar value representing the separability measure
56     """
57     classes = np.unique(y)
58     n_features = X.shape[1]
59
60     # Compute overall mean
61     overall_mean = np.mean(X, axis=0)
62
63     # Initialize scatter matrices
64     S_W = np.zeros((n_features, n_features)) # Within-class scatter
65     S_B = np.zeros((n_features, n_features)) # Between-class scatter
66
67     for c in classes:
68         X_c = X[y == c]
69         n_c = X_c.shape[0]
70
71         # Class mean
72         mean_c = np.mean(X_c, axis=0)
73
74         # Within-class scatter for class c
75         S_W += (X_c - mean_c).T @ (X_c - mean_c)
76
77         # Between-class scatter for class c
78         mean_diff = (mean_c - overall_mean).reshape(-1, 1)
79         S_B += n_c * (mean_diff @ mean_diff.T)
80
81     # Compute separability as trace( $S_B @ \text{inv}(S_W)$ )
82     # Add regularization to avoid singular matrix
83     regularization = 1e-6
84     S_W_reg = S_W + regularization * np.eye(n_features)
85
86     try:
87         separability = np.trace(np.linalg.solve(S_W_reg, S_B))
88     except np.linalg.LinAlgError:
89         logger.warning(
90             "Could not compute separability due to singular matrix. Returning 0."
91         )
92         separability = 0.0
93
94     return separability
95
96
97 def split_data(X: np.ndarray, y: np.ndarray, test_size=0.2, random_state=42):
98     """
99     Split the dataset into training and testing sets with stratification and feature scaling.
100     """
101
102     # Split the dataset
103     X_train, X_test, y_train, y_test = train_test_split(
104         X, y, test_size=test_size, random_state=random_state
105     )
106
107     # Feature Scaling
108     scaler = StandardScaler()
109     X_train = scaler.fit_transform(X_train)

```

```

110     X_test = scaler.transform(X_test)
111     logger.success("Feature scaling completed.\n")
112
113     return X_train, X_test, y_train, y_test
114
115
116 def baseline(X: np.ndarray, y: np.ndarray, dataset: str):
117     # split the dataset
118     X_train, X_test, y_train, y_test = split_data(
119         X, y, test_size=0.2, random_state=42
120     )
121
122     separability_before = compute_separability(X_train, y_train)
123
124     # Apply LDA
125     lda = LDA(n_components=None)
126     lda.fit(X_train, y_train)
127     X_train_lda = lda.transform(X_train)
128     X_test_lda = lda.transform(X_test)
129
130     separability_after = compute_separability(X_train_lda, y_train)
131
132     y_pred, y_score = lda.predict(X_test)
133     evaluate_classifier(
134         y_true=y_test,
135         y_pred=y_pred,
136         y_scores=y_score,
137         dataset_name=dataset,
138         classifier_name="LDA",
139         results_dir="results/task1/baseline",
140     )
141
142     logger.info(f"--- Task 1 ({dataset}) completed ---")
143     logger.info("Summary:")
144     logger.info(f"Separability before projection: {separability_before:.4f}")
145     logger.info(f"Separability after projection: {separability_after:.4f}")
146     logger.info("ROC/AUC plots saved to 'results/task1/baseline' directory.")
147
148
149 def dimension_sweep(X: np.ndarray, y: np.ndarray, dataset: str):
150     """
151     Sweep LDA projection dimensions d = 1, ..., C-1 and
152     record separability and test accuracy for each d.
153     Results are saved to results/task1/dim_sweep/{dataset}_dim_sweep.csv
154     """
155
156     # split the dataset
157     X_train, X_test, y_train, y_test = split_data(
158         X, y, test_size=0.2, random_state=42
159     )
160
161     classes = np.unique(y_train)
162     C = len(classes)
163     max_dim = min(C - 1, X_train.shape[1])
164
165     if max_dim < 1:
166         logger.error(
167             f"dim_sweep not applicable: C={C}, features={X_train.shape[1]}."
168         )
169         return
170
171     logger.info(
172         f"[{dataset}] dim_sweep: number of classes = {C}, "
173         f"feature dim = {X_train.shape[1]}, max_dim = {max_dim}"
174     )

```

```

175
176     results = []
177
178     for d in range(1, max_dim + 1):
179         lda = LDA(n_components=d)
180         lda.fit(X_train, y_train)
181
182         # Project to LDA subspace for computing separability
183         X_train_lda = lda.transform(X_train)
184         sep_d = compute_separability(X_train_lda, y_train)
185
186         # Classification: use LDA's decision rule in the original (scaled) space
187         y_pred, y_score = lda.predict(X_test)
188         acc_d = accuracy_score(y_test, y_pred)
189
190         logger.info(
191             f"[{dataset}] dim={d}: separability={sep_d:.4f}, "
192             f"accuracy={acc_d:.4f}"
193         )
194
195         results.append(
196             {
197                 "dim": d,
198                 "separability": sep_d,
199                 "accuracy": acc_d,
200             }
201         )
202
203     import os
204
205     os.makedirs("results/task1/dim_sweep", exist_ok=True)
206     df = pd.DataFrame(results)
207     out_path = f"results/task1/dim_sweep/{dataset}_dim_sweep.csv"
208     df.to_csv(out_path, index=False)
209
210     logger.success(
211         f"[{dataset}] dim_sweep finished. Results saved to '{out_path}'."
212     )
213
214
215 def prior_experiment(X: np.ndarray, y: np.ndarray, dataset: str):
216     """
217     Compare LDA with empirical priors vs balanced priors on an imbalanced dataset
218     (here, Taiwanese Bankruptcy).
219
220     Results (accuracy, balanced accuracy, precision, recall, F1, AUC) are saved to
221     results/task1/prior/{dataset}_prior.json
222     """
223
224     classes, counts = np.unique(y, return_counts=True)
225     C = len(classes)
226     if C != 2:
227         logger.error(
228             f"prior_experiment is designed for binary datasets, got C={C}."
229         )
230         return
231
232     logger.info(
233         f"[{dataset}] prior_experiment: classes={classes}, counts={counts}"
234     )
235
236     X_train, X_test, y_train, y_test = split_data(
237         X, y, test_size=0.2, random_state=42
238     )
239

```

```

240     # 1. empirical prior
241     lda_emp = LDA(n_components=None, priors=None)
242     lda_emp.fit(X_train, y_train)
243     y_pred_emp, y_score_emp = lda_emp.predict(X_test)
244
245     # 2. balanced prior
246     balanced_prior = np.ones(C) / C
247     lda_bal = LDA(n_components=None, priors=balanced_prior)
248     lda_bal.fit(X_train, y_train)
249     y_pred_bal, y_score_bal = lda_bal.predict(X_test)
250
251     metrics = {}
252
253     metrics["empirical"] = {
254         "accuracy": float(accuracy_score(y_test, y_pred_emp)),
255         "balanced_accuracy": float(balanced_accuracy_score(y_test, y_pred_emp)),
256         "precision": float(precision_score(y_test, y_pred_emp)),
257         "recall": float(recall_score(y_test, y_pred_emp)),
258         "f1": float(f1_score(y_test, y_pred_emp)),
259         "auc": float(roc_auc_score(y_test, y_score_emp[:, 1])),
260     }
261
262     metrics["balanced"] = {
263         "accuracy": float(accuracy_score(y_test, y_pred_bal)),
264         "balanced_accuracy": float(balanced_accuracy_score(y_test, y_pred_bal)),
265         "precision": float(precision_score(y_test, y_pred_bal)),
266         "recall": float(recall_score(y_test, y_pred_bal)),
267         "f1": float(f1_score(y_test, y_pred_bal)),
268         "auc": float(roc_auc_score(y_test, y_score_bal[:, 1])),
269     }
270
271     # save results
272     os.makedirs("results/task1/prior", exist_ok=True)
273     out_path = f"results/task1/prior/{dataset}_prior.json"
274
275     import json
276
277     with open(out_path, "w") as f:
278         json.dump(metrics, f, indent=2)
279
280     logger.info(f"[{dataset}] prior_experiment metrics:\n{metrics}")
281     logger.success(
282         f"[{dataset}] prior_experiment finished. Results saved to '{out_path}'."
283     )
284
285
286 if __name__ == "__main__":
287     parser = argparse.ArgumentParser(
288         description="Test Classifier with dimension reduction techniques"
289     )
290     parser.add_argument(
291         "--dataset",
292         type=str,
293         required=True,
294         choices=["telescope", "bankruptcy", "cover_type", "dry_bean"],
295         help="Dataset to load (e.g., telescope, bankruptcy, cover_type, dry_bean)",
296     )
297     parser.add_argument(
298         "--mode",
299         type=str,
300         required=True,
301         choices=["baseline", "dim_sweep", "prior"],
302         help="Experiment mode: baseline or dimension sweep or prior",
303     )
304

```

```

305     args = parser.parse_args()
306
307     # Load dataset
308     X, y = load_dataset(args.dataset)
309     logger.info(
310         f"Loaded {args.dataset} dataset: X shape={X.shape}, y shape={y.shape}\n"
311     )
312
313     # run the selected experiment mode
314     if args.mode == "baseline":
315         baseline(X, y, args.dataset)
316
317     elif args.mode == "dim_sweep":
318         dimension_sweep(X, y, args.dataset)
319
320     elif args.mode == "prior":
321         if args.dataset != "bankruptcy":
322             logger.error(
323                 "Prior mode is only applicable to the bankruptcy dataset."
324             )
325             exit(1)
326
327     prior_experiment(X, y, args.dataset)

```

task2.py

```

1  import argparse
2  from pathlib import Path
3  from typing import Tuple, Literal, Dict
4
5  import numpy as np
6  from loguru import logger
7  from sklearn.model_selection import train_test_split
8  from sklearn.preprocessing import StandardScaler
9  from sklearn.decomposition import PCA
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.svm import SVC
12 from sklearn.metrics import accuracy_score, balanced_accuracy_score, f1_score
13
14 from data_loader import load_dataset
15
16
17 DatasetName = Literal["bankruptcy", "dry_bean"]
18
19
20 # Utils
21 def parse_args() -> argparse.Namespace:
22     parser = argparse.ArgumentParser(
23         description="PR HW2 - Task 2: PCA experiments"
24     )
25
26     parser.add_argument(
27         "--dataset",
28         type=str,
29         choices=["bankruptcy", "dry_bean"],
30         required=True,
31         help="Which dataset to run on.",
32     )
33     parser.add_argument(
34         "--experiment",
35         type=str,
36         choices=[
37             "explained_variance", # Experiment 1

```

```

38         "dim_sweep", # Experiment 2
39         "clf_compare", # Experiment 3
40     ],
41     required=True,
42     help="Which experiment to run.",
43 )
44 parser.add_argument(
45     "--output-dir",
46     type=str,
47     default="results/task2",
48     help="Directory to save CSVs / logs for the report.",
49 )
50 parser.add_argument(
51     "--random-state",
52     type=int,
53     default=42,
54     help="Random seed for reproducibility.",
55 )
56 parser.add_argument(
57     "--test-size",
58     type=float,
59     default=0.2,
60     help="Test set ratio (for train/val/test splitting).",
61 )
62 parser.add_argument(
63     "--val-size",
64     type=float,
65     default=0.25,
66     help="Validation ratio within train+val "
67     "(e.g. 0.25 means train:val:test = 0.6:0.2:0.2).",
68 )
69 return parser.parse_args()
70
71
72 def ensure_output_dir(path_str: str) -> Path:
73     out_dir = Path(path_str)
74     out_dir.mkdir(parents=True, exist_ok=True)
75     return out_dir
76
77
78 def split_and_scale(
79     X: np.ndarray,
80     y: np.ndarray,
81     test_size: float,
82     val_size: float,
83     random_state: int,
84 ) -> Tuple[
85     np.ndarray, np.ndarray, np.ndarray, np.ndarray, np.ndarray, np.ndarray
86 ]:
87     """
88     Split into train/val/test and apply StandardScaler fitted on train only.
89     """
90     X_trainval, X_test, y_trainval, y_test = train_test_split(
91         X,
92         y,
93         test_size=test_size,
94         random_state=random_state,
95         stratify=y,
96     )
97     X_train, X_val, y_train, y_val = train_test_split(
98         X_trainval,
99         y_trainval,
100         test_size=val_size,
101         random_state=random_state,
102         stratify=y_trainval,

```

```

103     )
104
105     scaler = StandardScaler()
106     X_train = scaler.fit_transform(X_train)
107     X_val = scaler.transform(X_val)
108     X_test = scaler.transform(X_test)
109
110     return X_train, X_val, X_test, y_train, y_val, y_test
111
112
113 def build_logreg(random_state: int) -> LogisticRegression:
114     return LogisticRegression(
115         max_iter=1000,
116         random_state=random_state,
117         n_jobs=-1,
118     )
119
120
121 def build_rbf_svm(random_state: int) -> SVC:
122     return SVC(
123         kernel="rbf",
124         probability=False,
125         random_state=random_state,
126     )
127
128
129 def evaluate_clf(
130     clf,
131     X_tr: np.ndarray,
132     y_tr: np.ndarray,
133     X_val: np.ndarray,
134     y_val: np.ndarray,
135 ) -> Dict[str, float]:
136     clf.fit(X_tr, y_tr)
137     y_pred = clf.predict(X_val)
138
139     return {
140         "acc": accuracy_score(y_val, y_pred),
141         "balanced_acc": balanced_accuracy_score(y_val, y_pred),
142         "f1_macro": f1_score(y_val, y_pred, average="macro"),
143     }
144
145
146 # Experiment 1:
147 # Explained variance curve
148 def run_experiment_explained_variance(
149     dataset: DatasetName,
150     args: argparse.Namespace,
151 ) -> None:
152     logger.info(f"[Experiment 1] Explained variance on dataset = {dataset}")
153
154     X, y = load_dataset(dataset)
155     logger.info(
156         f"Loaded dataset {dataset}: X shape={X.shape}, y shape={y.shape}"
157     )
158
159     scaler = StandardScaler()
160     X_scaled = scaler.fit_transform(X)
161
162     n_features = X_scaled.shape[1]
163     pca = PCA(n_components=n_features)
164     pca.fit(X_scaled)
165
166     explained = pca.explained_variance_ratio_
167     cumulative = np.cumsum(explained)

```

```

168
169 out_dir = ensure_output_dir(args.output_dir)
170 out_csv = out_dir / f"explained_variance_{dataset}.csv"
171
172 import csv
173
174 with out_csv.open("w", newline="") as f:
175     writer = csv.writer(f)
176     writer.writerow(["k", "explained_var_ratio", "cumulative_var_ratio"])
177     for k, (e, c) in enumerate(zip(explained, cumulative), start=1):
178         writer.writerow([k, e, c])
179
180 logger.success(f"Saved explained variance curve to: {out_csv}")
181
182
183 # Experiment 2:
184 # Dimension sweep on Bankruptcy + LogReg
185 def run_experiment_dim_sweep(
186     dataset: DatasetName,
187     args: argparse.Namespace,
188 ) -> None:
189     if dataset != "bankruptcy":
190         logger.warning(
191             "Experiment 2 (dim_sweep) is mainly designed for the "
192             "Taiwanese Bankruptcy dataset. You are running it on "
193             f"{dataset}, which is fine but less emphasized in the report."
194         )
195
196 logger.info(f"[Experiment 2] PCA dimension sweep on dataset = {dataset}")
197
198 X, y = load_dataset(dataset)
199 logger.info(
200     f"Loaded dataset {dataset}: X shape={X.shape}, y shape={y.shape}"
201 )
202
203 X_train, X_val, X_test, y_train, y_val, y_test = split_and_scale(
204     X,
205     y,
206     test_size=args.test_size,
207     val_size=args.val_size,
208     random_state=args.random_state,
209 )
210
211 n_features = X_train.shape[1]
212
213 candidate_dims = [2, 5, 10, 20, 30, 40, 60, 80]
214 candidate_dims = [k for k in candidate_dims if k <= n_features]
215
216 out_dir = ensure_output_dir(args.output_dir)
217 out_csv = out_dir / f"dim_sweep_{dataset}.csv"
218
219 import csv
220
221 with out_csv.open("w", newline="") as f:
222     writer = csv.writer(f)
223     writer.writerow(
224         [
225             "k",
226             "val_acc",
227             "val_balanced_acc",
228             "val_f1_macro",
229             "test_acc",
230             "test_balanced_acc",
231             "test_f1_macro",
232         ]

```



```

233     )
234
235     # Baseline: no PCA
236     logger.info("Running baseline (no PCA)...")
237     clf_base = build_logreg(args.random_state)
238     metrics_val = evaluate_clf(clf_base, X_train, y_train, X_val, y_val)
239     metrics_test = evaluate_clf(
240         clf_base,
241         np.vstack([X_train, X_val]),
242         np.hstack([y_train, y_val]),
243         X_test,
244         y_test,
245     )
246
247     writer.writerow(
248         [
249             0,
250             metrics_val["acc"],
251             metrics_val["balanced_acc"],
252             metrics_val["f1_macro"],
253             metrics_test["acc"],
254             metrics_test["balanced_acc"],
255             metrics_test["f1_macro"],
256         ]
257     )
258
259     logger.success(
260         f"Baseline (no PCA) - "
261         f"val_acc={metrics_val['acc']:.4f}, "
262         f"val_balanced_acc={metrics_val['balanced_acc']:.4f}"
263     )
264
265     # PCA + LogReg for each k
266     for k in candidate_dims:
267         logger.info(f"Running PCA + LogReg with k={k}...")
268         pca = PCA(n_components=k, random_state=args.random_state)
269         X_train_pca = pca.fit_transform(X_train)
270         X_val_pca = pca.transform(X_val)
271         X_test_pca = pca.transform(X_test)
272
273         clf = build_logreg(args.random_state)
274         metrics_val = evaluate_clf(
275             clf, X_train_pca, y_train, X_val_pca, y_val
276         )
277         metrics_test = evaluate_clf(
278             clf,
279             np.vstack([X_train_pca, X_val_pca]),
280             np.hstack([y_train, y_val]),
281             X_test_pca,
282             y_test,
283         )
284
285         writer.writerow(
286             [
287                 k,
288                 metrics_val["acc"],
289                 metrics_val["balanced_acc"],
290                 metrics_val["f1_macro"],
291                 metrics_test["acc"],
292                 metrics_test["balanced_acc"],
293                 metrics_test["f1_macro"],
294             ]
295         )
296
297     logger.success(

```

```

298         f"k={k}: val_acc={metrics_val['acc']:.4f}, "
299         f"val_balanced_acc={metrics_val['balanced_acc']:.4f}"
300     )
301
302     logger.success(f"Saved dimension sweep results to: {out_csv}")
303
304
305     # Experiment 3:
306     # Compare PCA effects across classifiers & datasets
307     def run_experiment_classifier_compare(
308         dataset: DatasetName,
309         args: argparse.Namespace,
310     ) -> None:
311         logger.info(f"[Experiment 3] Classifier comparison on dataset = {dataset}")
312
313         X, y = load_dataset(dataset)
314         logger.info(
315             f"Loaded dataset {dataset}: X shape={X.shape}, y shape={y.shape}"
316         )
317
318         X_train, X_val, X_test, y_train, y_val, y_test = split_and_scale(
319             X,
320             y,
321             test_size=args.test_size,
322             val_size=args.val_size,
323             random_state=args.random_state,
324         )
325
326         out_dir = ensure_output_dir(args.output_dir)
327         out_csv = out_dir / f"classifier_compare_{dataset}.csv"
328
329         import csv
330
331         with out_csv.open("w", newline="") as f:
332             writer = csv.writer(f)
333             writer.writerow(
334                 [
335                     "classifier",
336                     "use_pca",
337                     "pca_dim",
338                     "val_acc",
339                     "val_balanced_acc",
340                     "val_f1_macro",
341                     "test_acc",
342                     "test_balanced_acc",
343                     "test_f1_macro",
344                 ]
345             )
346
347         # 1. Baseline: LogReg / RBF-SVM without PCA
348         for clf_name in ["logreg", "rbf_svm"]:
349             logger.info(f"Running baseline without PCA for {clf_name}...")
350
351             if clf_name == "logreg":
352                 clf = build_logreg(args.random_state)
353             else:
354                 clf = build_rbf_svm(args.random_state)
355
356             metrics_val = evaluate_clf(clf, X_train, y_train, X_val, y_val)
357             metrics_test = evaluate_clf(
358                 clf,
359                 np.vstack([X_train, X_val]),
360                 np.hstack([y_train, y_val]),
361                 X_test,
362                 y_test,

```

```

363         )
364
365     writer.writerow(
366         [
367             clf_name,
368             0,
369             0,
370             metrics_val["acc"],
371             metrics_val["balanced_acc"],
372             metrics_val["f1_macro"],
373             metrics_test["acc"],
374             metrics_test["balanced_acc"],
375             metrics_test["f1_macro"],
376         ]
377     )
378
379     # 2. PCA + classifier
380     logger.info("Fitting PCA (95% variance) for classifier comparison...")
381     pca = PCA(
382         n_components=0.95, svd_solver="full", random_state=args.random_state
383     )
384     X_train_pca = pca.fit_transform(X_train)
385     X_val_pca = pca.transform(X_val)
386     X_test_pca = pca.transform(X_test)
387     k_eff = X_train_pca.shape[1]
388
389     logger.info(f"Effective PCA dimension for 95% variance: k={k_eff}")
390
391     for clf_name in ["logreg", "rbf_svm"]:
392         logger.info(f"Running {clf_name} with PCA (k={k_eff})...")
393
394         if clf_name == "logreg":
395             clf = build_logreg(args.random_state)
396         else:
397             clf = build_rbf_svm(args.random_state)
398
399         metrics_val = evaluate_clf(
400             clf, X_train_pca, y_train, X_val_pca, y_val
401         )
402         metrics_test = evaluate_clf(
403             clf,
404             np.vstack([X_train_pca, X_val_pca]),
405             np.hstack([y_train, y_val]),
406             X_test_pca,
407             y_test,
408         )
409
410         writer.writerow(
411             [
412                 clf_name,
413                 1,
414                 k_eff,
415                 metrics_val["acc"],
416                 metrics_val["balanced_acc"],
417                 metrics_val["f1_macro"],
418                 metrics_test["acc"],
419                 metrics_test["balanced_acc"],
420                 metrics_test["f1_macro"],
421             ]
422         )
423
424     logger.success(f"Saved classifier comparison results to: {out_csv}")
425
426
427 if __name__ == "__main__":

```

```

428     args = parse_args()
429
430     logger.info(
431         f"Running Task 2 with dataset={args.dataset}, "
432         f"experiment={args.experiment}"
433     )
434
435     if args.experiment == "explained_variance":
436         run_experiment_explained_variance(args.dataset, args)
437     elif args.experiment == "dim_sweep":
438         run_experiment_dim_sweep(args.dataset, args)
439     elif args.experiment == "clf_compare":
440         run_experiment_classifier_compare(args.dataset, args)
441     else:
442         raise ValueError(f"Unknown experiment: {args.experiment}")

```

task2_plot.py

```

1  import pandas as pd
2  import matplotlib.pyplot as plt
3
4  bank = pd.read_csv("results/task2/explained_variance_bankruptcy.csv")
5  bean = pd.read_csv("results/task2/explained_variance_dry_bean.csv")
6
7  plt.figure(figsize=(8,5))
8  plt.plot(bank["k"], bank["cumulative_var_ratio"], label="Bankruptcy")
9  plt.plot(bean["k"], bean["cumulative_var_ratio"], label="Dry Bean")
10
11  plt.xlabel("Number of Principal Components $k$")
12  plt.ylabel("Cumulative Explained Variance Ratio")
13  plt.title("Explained Variance Curve (PCA)")
14  plt.grid(True, linestyle="--", alpha=0.4)
15  plt.legend()
16
17  plt.tight_layout()
18  # plt.savefig("figs/task2_explained_variance.pdf")
19  plt.savefig("results/task2/explained_variance.png", dpi=300)

```

task3.py

```

1  #!/usr/bin/env python
2  import argparse
3  from pathlib import Path
4  from typing import Literal
5
6  import numpy as np
7  from loguru import logger
8  from sklearn.preprocessing import StandardScaler
9  from sklearn.decomposition import PCA
10  from sklearn.manifold import Isomap
11  from sklearn.model_selection import train_test_split
12  from sklearn.metrics import accuracy_score, balanced_accuracy_score, f1_score
13  from sklearn.linear_model import LogisticRegression
14  from sklearn.svm import SVC
15
16  from data_loader import load_dataset
17
18
19  DatasetName = Literal["bankruptcy", "dry_bean", "telescope", "cover_type"]
20
21

```

```

22 # Utils
23 def parse_args() -> argparse.Namespace:
24     parser = argparse.ArgumentParser(
25         description="PR HW2 - Task 3: Nonlinear dimensionality reduction (Isomap) experiments"
26     )
27
28     parser.add_argument(
29         "--dataset",
30         type=str,
31         choices=["bankruptcy", "dry_bean", "telescope", "cover_type"],
32         default="dry_bean",
33         help="Which dataset to run on (default: dry_bean).",
34     )
35     parser.add_argument(
36         "--experiment",
37         type=str,
38         choices=[
39             "embedding_viz", # Experiment 1: 2D embedding visualization
40             "dim_sweep_isomap", # Experiment 2: Isomap dimension sweep (to be implemented)
41             "neighbor_sensitivity", # Experiment 3: n_neighbors sensitivity (to be implemented)
42         ],
43         required=True,
44         help="Which experiment to run.",
45     )
46     parser.add_argument(
47         "--output-dir",
48         type=str,
49         default="results/task3",
50         help="Directory to save figures / CSVs.",
51     )
52     parser.add_argument(
53         "--random-state",
54         type=int,
55         default=42,
56         help="Random seed for reproducibility.",
57     )
58     parser.add_argument(
59         "--n-neighbors",
60         type=int,
61         default=10,
62         help="Number of neighbors for Isomap (default: 10).",
63     )
64     return parser.parse_args()
65
66
67 def ensure_output_dir(path_str: str) -> Path:
68     out_dir = Path(path_str)
69     out_dir.mkdir(parents=True, exist_ok=True)
70     return out_dir
71
72
73 def split_and_scale(
74     X: np.ndarray,
75     y: np.ndarray,
76     random_state: int,
77 ):
78     """
79     Split into train/val/test and apply StandardScaler fitted on train only.
80     We use the same 60/20/20 split protocol as in Task 2:
81     - First split: train+val vs test = 0.8 / 0.2
82     - Second split: train vs val = 0.75 / 0.25 of (train+val)
83     """
84     X_trainval, X_test, y_trainval, y_test = train_test_split(
85         X,
86         y,

```

```

87         test_size=0.2,
88         random_state=random_state,
89         stratify=y,
90     )
91
92     X_train, X_val, y_train, y_val = train_test_split(
93         X_trainval,
94         y_trainval,
95         test_size=0.25, # 0.8 * 0.25 = 0.2 => train:val:test = 0.6:0.2:0.2
96         random_state=random_state,
97         stratify=y_trainval,
98     )
99
100     scaler = StandardScaler()
101     X_train = scaler.fit_transform(X_train)
102     X_val = scaler.transform(X_val)
103     X_test = scaler.transform(X_test)
104
105     return X_train, X_val, X_test, y_train, y_val, y_test
106
107
108 def build_logreg(random_state: int) -> LogisticRegression:
109     return LogisticRegression(
110         max_iter=1000,
111         random_state=random_state,
112         n_jobs=-1,
113     )
114
115
116 def build_rbf_svm(random_state: int) -> SVC:
117     return SVC(
118         kernel="rbf",
119         probability=False,
120         random_state=random_state,
121     )
122
123
124 def evaluate_clf(
125     clf,
126     X_tr: np.ndarray,
127     y_tr: np.ndarray,
128     X_val: np.ndarray,
129     y_val: np.ndarray,
130 ):
131     clf.fit(X_tr, y_tr)
132     y_pred = clf.predict(X_val)
133
134     return {
135         "acc": accuracy_score(y_val, y_pred),
136         "balanced_acc": balanced_accuracy_score(y_val, y_pred),
137         "f1_macro": f1_score(y_val, y_pred, average="macro"),
138     }
139
140
141 # Experiment 1:
142 # 2D Embedding Visualization (PCA vs Isomap)
143 def run_experiment_embedding_viz(
144     dataset: DatasetName,
145     args: argparse.Namespace,
146 ) -> None:
147     """
148     Experiment 1:
149     Compare 2D embeddings produced by PCA and Isomap on a given dataset.
150     This is mainly designed for Dry Bean, but can be run on other datasets as well.
151     """

```

```

152     if dataset != "dry_bean":
153         logger.warning(
154             "Experiment 1 (embedding_viz) is primarily designed for the Dry Bean dataset, "
155             f"but you are running it on {dataset}."
156         )
157
158     logger.info(
159         f"[Experiment 1] 2D embedding visualization on dataset = {dataset}"
160     )
161
162     # Load full dataset (no split; this is for visualization)
163     X, y = load_dataset(dataset)
164     logger.info(
165         f"Loaded dataset {dataset}: X shape={X.shape}, y shape={y.shape}"
166     )
167
168     # Standardize features
169     scaler = StandardScaler()
170     X_scaled = scaler.fit_transform(X)
171
172     # 2D PCA
173     logger.info("Fitting PCA (2D)...")
174     pca = PCA(n_components=2, random_state=args.random_state)
175     X_pca = pca.fit_transform(X_scaled)
176
177     # 2D Isomap
178     logger.info(f"Fitting Isomap (2D) with n_neighbors={args.n_neighbors}...")
179     iso = Isomap(
180         n_components=2,
181         n_neighbors=args.n_neighbors,
182     )
183     X_iso = iso.fit_transform(X_scaled)
184
185     out_dir = ensure_output_dir(args.output_dir)
186
187     # Save the embeddings as CSV for later plotting or analysis
188     import csv
189
190     pca_csv = out_dir / f"embedding_{dataset}_pca2.csv"
191     iso_csv = out_dir / f"embedding_{dataset}_isomap2_k{args.n_neighbors}.csv"
192
193     logger.info(f"Saving PCA 2D embedding to: {pca_csv}")
194     with pca_csv.open("w", newline="") as f:
195         writer = csv.writer(f)
196         writer.writerow(["z1", "z2", "label"])
197         for (z1, z2), label in zip(X_pca, y):
198             writer.writerow([z1, z2, label])
199
200     logger.info(f"Saving Isomap 2D embedding to: {iso_csv}")
201     with iso_csv.open("w", newline="") as f:
202         writer = csv.writer(f)
203         writer.writerow(["z1", "z2", "label"])
204         for (z1, z2), label in zip(X_iso, y):
205             writer.writerow([z1, z2, label])
206
207     # Optionally, generate a side-by-side scatter plot for report
208     try:
209         import matplotlib.pyplot as plt
210
211         logger.info("Generating 2D scatter plots for PCA and Isomap...")
212
213         classes = np.unique(y)
214
215         fig, axes = plt.subplots(1, 2, figsize=(10, 4))
216

```

```

217     # PCA plot
218     ax = axes[0]
219     for c in classes:
220         idx = y == c
221         ax.scatter(
222             X_pca[idx, 0],
223             X_pca[idx, 1],
224             s=5,
225             alpha=0.6,
226             label=str(c),
227         )
228     ax.set_title("PCA (2D)")
229     ax.set_xlabel("PC1")
230     ax.set_ylabel("PC2")
231
232     # Isomap plot
233     ax = axes[1]
234     for c in classes:
235         idx = y == c
236         ax.scatter(
237             X_iso[idx, 0],
238             X_iso[idx, 1],
239             s=5,
240             alpha=0.6,
241             label=str(c),
242         )
243     ax.set_title(f"Isomap (2D, k={args.n_neighbors})")
244     ax.set_xlabel("Component 1")
245     ax.set_ylabel("Component 2")
246
247     # Legend: 只在右圖畫一次，避免擠爆
248     handles, labels = axes[1].get_legend_handles_labels()
249     fig.legend(
250         handles,
251         labels,
252         loc="upper center",
253         ncol=min(len(classes), 6),
254         fontsize=8,
255     )
256
257     plt.tight_layout(rect=[0, 0, 1, 0.90])
258
259     fig_path_pdf = out_dir / f"task3_{dataset}_pca_isomap_2d.pdf"
260     fig_path_png = out_dir / f"task3_{dataset}_pca_isomap_2d.png"
261
262     fig.savefig(fig_path_pdf)
263     fig.savefig(fig_path_png, dpi=300)
264
265     plt.close(fig)
266
267     logger.success(
268         f"Saved 2D embedding figure to: {fig_path_pdf} and {fig_path_png}"
269     )
270
271     except ImportError:
272         logger.warning(
273             "matplotlib is not available; skipped figure generation. "
274             "You can still use the saved CSV embeddings for plotting."
275         )
276
277
278     # Experiment 2:
279     # Isomap Dimension Sweep
280     def run_experiment_dim_sweep_isomap(
281         dataset: DatasetName,

```



```

282     args: argparse.Namespace,
283 ) -> None:
284     """
285     Isomap dimension sweep experiment.
286
287     - Split dataset into train/val/test using the same 60/20/20 protocol as Task 2.
288     - For a range of Isomap dimensions (e.g., k {2, 5, 10, 20, 30}),
289       fit Isomap on the training data (after standardization) and
290       evaluate logistic regression on val/test.
291     - Save all metrics to CSV.
292     """
293     if dataset != "bankruptcy":
294         logger.warning(
295             "Experiment 'dim_sweep_isomap' is mainly designed for the "
296             "Taiwanese Bankruptcy dataset, but you are running it on "
297             f"{dataset}."
298         )
299
300     logger.info(f"[Experiment 2] Isomap dimension sweep on dataset = {dataset}")
301
302     # 1. Load dataset
303     X, y = load_dataset(dataset)
304     logger.info(f"Loaded dataset {dataset}: X shape={X.shape}, y shape={y.shape}")
305
306     # 2. Split + scale
307     X_train, X_val, X_test, y_train, y_val, y_test = split_and_scale(
308         X,
309         y,
310         random_state=args.random_state,
311     )
312
313     n_features = X_train.shape[1]
314
315     # 3. Candidate Isomap dimensions (bounded by #features)
316     candidate_dims = [2, 5, 10, 20, 30]
317     candidate_dims = [k for k in candidate_dims if k <= n_features]
318
319     out_dir = ensure_output_dir(args.output_dir)
320     out_csv = out_dir / f"dim_sweep_isomap_{dataset}.csv"
321
322     import csv
323
324     with out_csv.open("w", newline="") as f:
325         writer = csv.writer(f)
326         writer.writerow(
327             [
328                 "k",
329                 "val_acc",
330                 "val_balanced_acc",
331                 "val_f1_macro",
332                 "test_acc",
333                 "test_balanced_acc",
334                 "test_f1_macro",
335             ]
336         )
337
338     # 4. Baseline: no dimensionality reduction (just logistic regression on scaled space)
339     logger.info("Running baseline (no dimensionality reduction)...")
340     clf_base = build_logreg(args.random_state)
341     metrics_val = evaluate_clf(clf_base, X_train, y_train, X_val, y_val)
342     metrics_test = evaluate_clf(
343         clf_base,
344         np.vstack([X_train, X_val]),
345         np.hstack([y_train, y_val]),
346         X_test,

```

```

347         y_test,
348     )
349
350     writer.writerow(
351         [
352             0,
353             metrics_val["acc"],
354             metrics_val["balanced_acc"],
355             metrics_val["f1_macro"],
356             metrics_test["acc"],
357             metrics_test["balanced_acc"],
358             metrics_test["f1_macro"],
359         ]
360     )
361
362     logger.success(
363         "Baseline (no DR) - "
364         f"val_bal_acc={metrics_val['balanced_acc']:.4f}, "
365         f"test_bal_acc={metrics_test['balanced_acc']:.4f}"
366     )
367
368     # 5. For each Isomap dimension, fit Isomap + LogReg
369     for k in candidate_dims:
370         logger.info(
371             f"Running Isomap + Logistic Regression with n_components={k}, "
372             f"n_neighbors={args.n_neighbors}..."
373         )
374
375         iso = Isomap(
376             n_components=k,
377             n_neighbors=args.n_neighbors,
378         )
379
380         # Fit Isomap on training data only, then transform val/test
381         X_train_iso = iso.fit_transform(X_train)
382         X_val_iso = iso.transform(X_val)
383         X_test_iso = iso.transform(X_test)
384
385         clf = build_logreg(args.random_state)
386
387         metrics_val = evaluate_clf(
388             clf,
389             X_train_iso,
390             y_train,
391             X_val_iso,
392             y_val,
393         )
394         metrics_test = evaluate_clf(
395             clf,
396             np.vstack([X_train_iso, X_val_iso]),
397             np.hstack([y_train, y_val]),
398             X_test_iso,
399             y_test,
400         )
401
402         writer.writerow(
403             [
404                 k,
405                 metrics_val["acc"],
406                 metrics_val["balanced_acc"],
407                 metrics_val["f1_macro"],
408                 metrics_test["acc"],
409                 metrics_test["balanced_acc"],
410                 metrics_test["f1_macro"],
411             ]

```

```

412         )
413
414         logger.success(
415             f"k={k}: val_bal_acc={metrics_val['balanced_acc']:.4f}, "
416             f"test_bal_acc={metrics_test['balanced_acc']:.4f}"
417         )
418
419     logger.success(f"Saved Isomap dimension sweep results to: {out_csv}")
420
421
422     # Experiment 3:
423     def run_experiment_neighbor_sensitivity(
424         dataset: DatasetName,
425         args: argparse.Namespace,
426     ) -> None:
427         """
428         Hyperparameter sensitivity experiment for Isomap.
429
430         - Fix Isomap output dimension (n_components = 10).
431         - Vary n_neighbors (e.g., {5, 10, 20, 50}).
432         - Evaluate performance of an RBF SVM classifier on val/test for each setting.
433         - Also include a baseline RBF SVM without Isomap.
434         - Save all metrics to CSV.
435         """
436         logger.info(
437             f"[Experiment 3] Isomap neighbor sensitivity on dataset = {dataset}"
438         )
439
440         # 1. Load dataset
441         X, y = load_dataset(dataset)
442         logger.info(f"Loaded dataset {dataset}: X shape={X.shape}, y shape={y.shape}")
443
444         # 2. Split + scale (same 60/20/20 protocol)
445         X_train, X_val, X_test, y_train, y_val, y_test = split_and_scale(
446             X,
447             y,
448             random_state=args.random_state,
449         )
450
451         out_dir = ensure_output_dir(args.output_dir)
452         out_csv = out_dir / f"neighbor_sensitivity_isomap_{dataset}.csv"
453
454         # Choose a reasonable fixed output dimension for Isomap
455         n_components = min(10, X_train.shape[1])
456
457         # Candidate neighbor sizes
458         candidate_neighbors = [5, 10, 20, 50]
459
460         import csv
461
462         with out_csv.open("w", newline="") as f:
463             writer = csv.writer(f)
464             writer.writerow(
465                 [
466                     "use_isomap",
467                     "n_components",
468                     "n_neighbors",
469                     "val_acc",
470                     "val_balanced_acc",
471                     "val_f1_macro",
472                     "test_acc",
473                     "test_balanced_acc",
474                     "test_f1_macro",
475                 ]
476             )

```

```

477
478 # 3. Baseline: RBF SVM without Isomap
479 logger.info("Running baseline RBF SVM (no Isomap)...")
480 clf_base = build_rbf_svm(args.random_state)
481 metrics_val = evaluate_clf(clf_base, X_train, y_train, X_val, y_val)
482 metrics_test = evaluate_clf(
483     clf_base,
484     np.vstack([X_train, X_val]),
485     np.hstack([y_train, y_val]),
486     X_test,
487     y_test,
488 )
489
490 writer.writerow(
491     [
492         0,          # use_isomap = 0
493         0,          # n_components
494         0,          # n_neighbors
495         metrics_val["acc"],
496         metrics_val["balanced_acc"],
497         metrics_val["f1_macro"],
498         metrics_test["acc"],
499         metrics_test["balanced_acc"],
500         metrics_test["f1_macro"],
501     ]
502 )
503
504 logger.success(
505     "Baseline RBF SVM - "
506     f"val_bal_acc={metrics_val['balanced_acc']:.4f}, "
507     f"test_bal_acc={metrics_test['balanced_acc']:.4f}"
508 )
509
510 # 4. Isomap + RBF SVM for each neighbor setting
511 for k_nn in candidate_neighbors:
512     logger.info(
513         f"Running Isomap (n_components={n_components}, "
514         f"n_neighbors={k_nn}) + RBF SVM..."
515     )
516
517     iso = Isomap(
518         n_components=n_components,
519         n_neighbors=k_nn,
520     )
521
522     # Fit Isomap on training data only
523     X_train_iso = iso.fit_transform(X_train)
524     X_val_iso = iso.transform(X_val)
525     X_test_iso = iso.transform(X_test)
526
527     clf = build_rbf_svm(args.random_state)
528
529     metrics_val = evaluate_clf(
530         clf,
531         X_train_iso,
532         y_train,
533         X_val_iso,
534         y_val,
535     )
536     metrics_test = evaluate_clf(
537         clf,
538         np.vstack([X_train_iso, X_val_iso]),
539         np.hstack([y_train, y_val]),
540         X_test_iso,
541         y_test,

```

```

542         )
543
544     writer.writerow(
545         [
546             1,          # use_isomap = 1
547             n_components, # n_components
548             k_nn,        # n_neighbors
549             metrics_val["acc"],
550             metrics_val["balanced_acc"],
551             metrics_val["f1_macro"],
552             metrics_test["acc"],
553             metrics_test["balanced_acc"],
554             metrics_test["f1_macro"],
555         ]
556     )
557
558     logger.success(
559         f"n_neighbors={k_nn}: "
560         f"val_bal_acc={metrics_val['balanced_acc']:.4f}, "
561         f"test_bal_acc={metrics_test['balanced_acc']:.4f}"
562     )
563
564     logger.success(
565         f"Saved Isomap neighbor sensitivity results to: {out_csv}"
566     )
567
568
569 if __name__ == "__main__":
570     args = parse_args()
571
572     logger.info(
573         f"Running Task 3 with dataset={args.dataset}, "
574         f"experiment={args.experiment}, "
575         f"n_neighbors={args.n_neighbors}"
576     )
577
578     if args.experiment == "embedding_viz":
579         run_experiment_embedding_viz(args.dataset, args)
580     elif args.experiment == "dim_sweep_isomap":
581         run_experiment_dim_sweep_isomap(args.dataset, args)
582     elif args.experiment == "neighbor_sensitivity":
583         run_experiment_neighbor_sensitivity(args.dataset, args)
584     else:
585         raise ValueError(f"Unknown experiment: {args.experiment}")

```