

NYCU Pattern Recognition, Assignment #1

Student ID: 112550157, Name: 黃冠瑋

1 Methodology

Our implementation consists of several Python modules designed to fulfill the assignment requirements and enable a flexible experimental workflow.

1.1 Data Loading (`data_loader.py`)

We use the `ucimlrepo` package to retrieve the four main datasets from the UCI Machine Learning Repository: telescope, bankruptcy, dry bean, and cover type. To accelerate repeated data access, a caching mechanism is implemented via `joblib.Memory`, which stores preprocessed datasets locally.

The core function `_preprocess_data()` standardizes features and encodes categorical labels using `sklearn.preprocessing.LabelEncoder`. This ensures that all datasets are numerical and compatible with the implemented classifiers.

1.2 Classifier Module (`classifier.py`)

Following the assignment specification, each classifier is implemented as an independent Python class, including MLP, Naive Bayes, SVM, KNN, Random Forest, and XGBoost.

Each class defines two key methods:

- *train*(X, y) —fits the model using the given training data.
- *predict*(X) - outputs both the predicted class labels (y_{pred}) and discriminant function values (y_{scores} , e.g., class probabilities).

1.3 Evaluation Module (evaluation.py)

The evaluation process is implemented in the function `evaluate_classifier()`, which performs both quantitative and visual analyses.

- **Confusion Matrix:** Computed via `sklearn.metrics.confusion_matrix` and visualized as a heatmap using Seaborn, then saved as `.png`.
- **Classification Report:** Summarizes precision, recall, F1-score, and accuracy via `sklearn.metrics.classification_report`.
- **ROC Curve and AUC:** For binary classification tasks, the function automatically generates the ROC curve and computes the AUC (Area Under Curve), also saved as `.png`.

These metrics collectively assess both overall accuracy and class-level performance.

1.4 Experimental Workflow (main.py)

The main script coordinates dataset loading, model training, and evaluation. Before cross-validation, the dataset is divided into training and testing subsets using an 80/20 stratified split to preserve class balance.

The workflow then proceeds in two stages:

- **Cross Validation:** The function `k_fold_cross_validation()` performs 5-fold cross-validation to estimate model stability and prevent overfitting.
- **Final Evaluation:** After cross-validation, each model is retrained on the full dataset and re-evaluated on the same data to report its final performance metrics.

2 Experiments and Results

2.1 Experimental Setup

All experiments are conducted using Python 3.12 with scikit-learn and XGBoost. Each dataset is normalized and split with a 80/20 train-test ratio while maintaining class balance. Performance is measured by accuracy, confusion matrix, and ROC—AUC (for binary datasets). Each result represents the mean over 5-fold cross-validation.

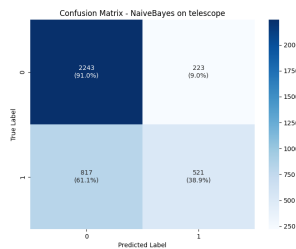
2.2 Dataset-wise Performance

2.2.1 MAGIC Gamma Telescope

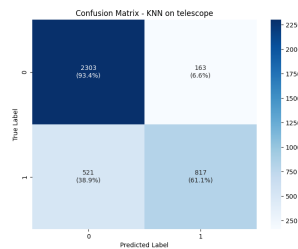
The MAGIC Gamma Telescope dataset is a binary classification task, used to classify high-energy gamma particles. It is a multivariate dataset containing 19020 instances, each described by 10 real-valued features. The performance of the classifiers on this dataset is summarized in Table 1.

Table 1: Telescope dataset: classification performance summary.

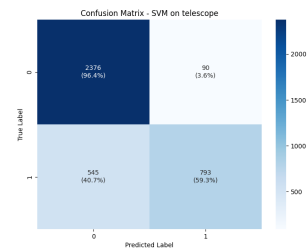
Classifier	CV Accuracy	Accuracy (full)	AUC
Naive Bayes	0.7257 ± 0.0103	0.73	0.7593
KNN	0.8013 ± 0.0044	0.82	0.8441
SVM	0.8198 ± 0.0067	0.83	0.8718
MLP	0.8189 ± 0.0222	0.84	0.8992
Random Forest	0.8761 ± 0.0029	0.89	0.9378
XGBoost	0.8745 ± 0.0030	0.88	0.9382



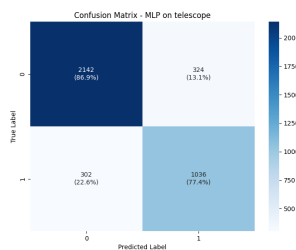
(a) Naive Bayes



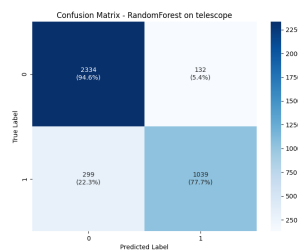
(b) KNN



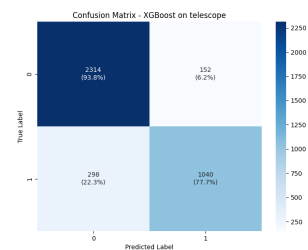
(c) SVM



(d) MLP

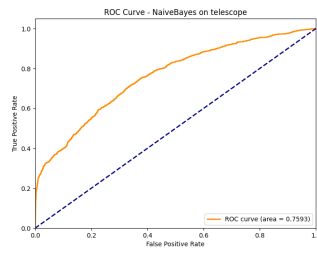


(e) Random Forest

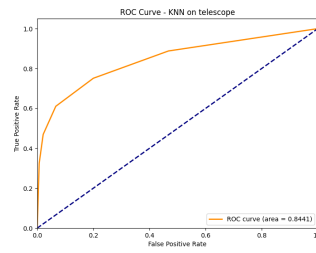


(f) XGBoost

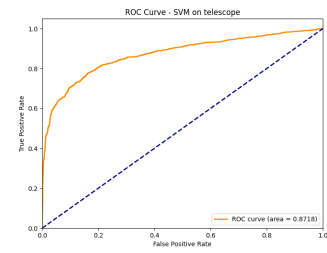
Figure 1: Confusion matrices for the Telescope dataset.



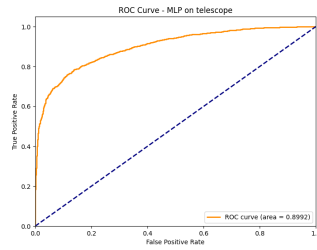
(a) Naive Bayes (AUC=0.7593)



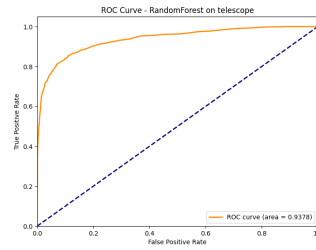
(b) KNN (AUC=0.8441)



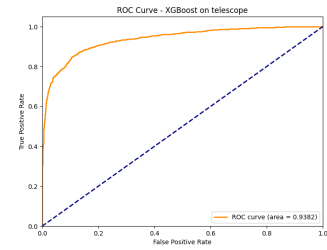
(c) SVM (AUC=0.8718)



(d) MLP (AUC=0.8992)



(e) Random Forest (AUC=0.9378)



(f) XGBoost (AUC=0.9382)

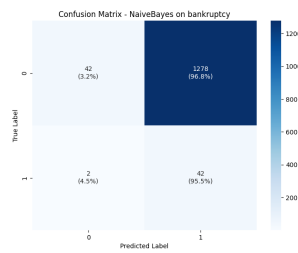
Figure 2: ROC curves and AUC scores for the Telescope dataset.

2.2.2 Taiwanese Bankruptcy Prediction

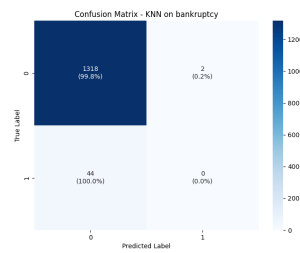
This dataset is a binary classification problem focused on business. It contains 6819 instances, each characterized by a high dimensionality of 95 integer features. The performance of the classifiers on this dataset, shown in Table 2, reveals significant discrepancies.

Table 2: Bankruptcy dataset: classification performance summary.

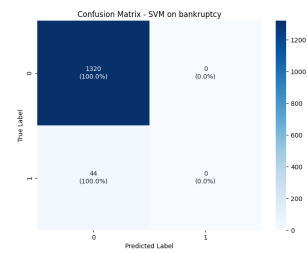
Classifier	CV Accuracy	Accuracy (full)	AUC
Naive Bayes	0.0631 ± 0.0069	0.06	0.6521
KNN	0.9665 ± 0.0012	0.97	0.5540
SVM	0.9677 ± 0.0004	0.97	0.5409
MLP	0.9419 ± 0.0171	0.94	0.5596
Random Forest	0.9696 ± 0.0023	0.97	0.9448
XGBoost	0.9696 ± 0.0051	0.97	0.9671



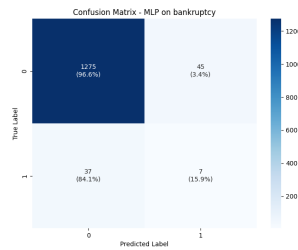
(a) Naive Bayes



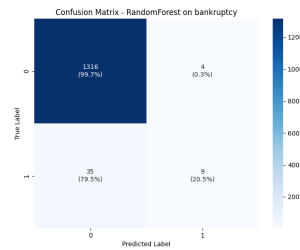
(b) KNN



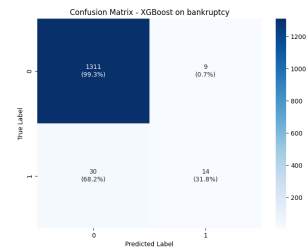
(c) SVM



(d) MLP

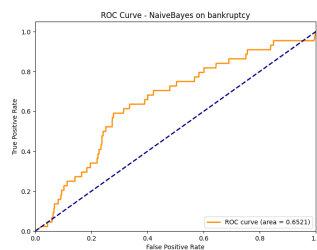


(e) Random Forest

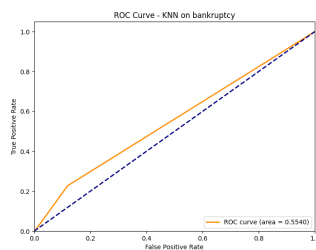


(f) XGBoost

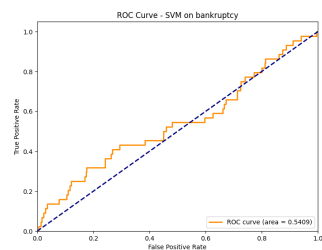
Figure 3: Confusion matrices for the bankruptcy dataset.



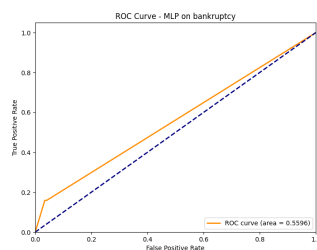
(a) Naive Bayes (AUC=0.7593)



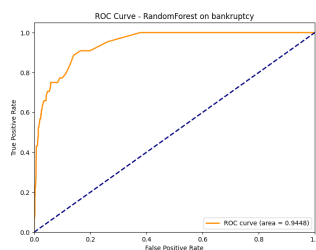
(b) KNN (AUC=0.8441)



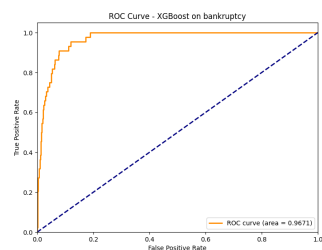
(c) SVM (AUC=0.8718)



(d) MLP (AUC=0.8992)



(e) Random Forest (AUC=0.9378)



(f) XGBoost (AUC=0.9382)

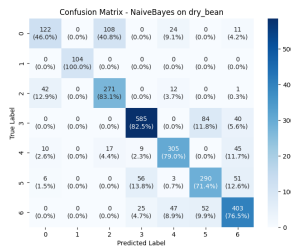
Figure 4: ROC curves and AUC scores for the bankruptcy dataset.

2.2.3 Dry Bean

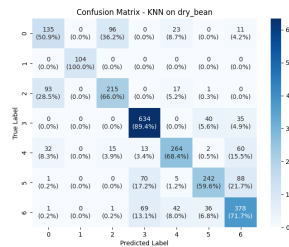
The Dry Bean dataset is a multi-class classification problem, with the goal of classifying 7 different types of registered dry beans. The dataset consists of 13611 instances, each described by 16 real-valued and integer features derived from images. Since this is a multi-class dataset, AUC is not a primary metric. The performance of the classifiers on this dataset is summarized in Table 3.

Table 3: Dry Bean dataset: classification performance summary.

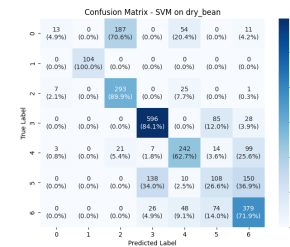
Classifier	CV Accuracy	Accuracy (full)
Naive Bayes	0.7639 ± 0.0055	0.76
KNN	0.7134 ± 0.0076	0.72
SVM	0.6336 ± 0.0019	0.64
MLP	0.4695 ± 0.1201	0.28
Random Forest	0.9232 ± 0.0054	0.92
XGBoost	0.9291 ± 0.0055	0.92



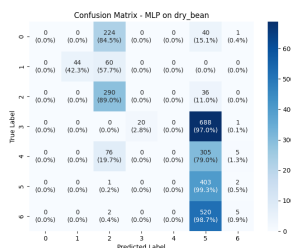
(a) Naive Bayes



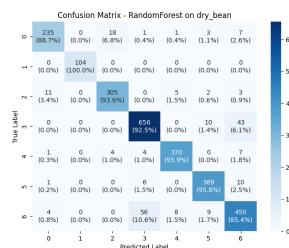
(b) KNN



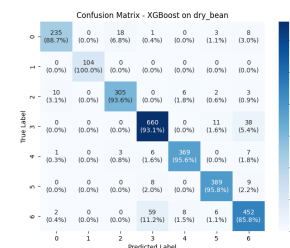
(c) SVM



(d) MLP



(e) Random Forest



(f) XGBoost

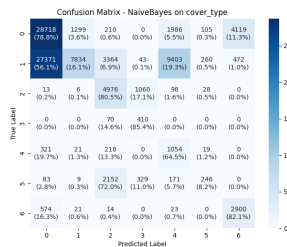
Figure 5: Confusion matrices for the dry bean dataset.

2.2.4 Cover Type

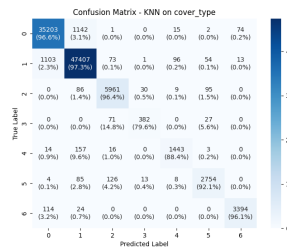
This is a large multi-class classification dataset with 7 distinct forest cover types. The full dataset has 581,012 instances and 54 features (both categorical and integer). Due to its size, we performed our experiments on a random subsample of 50000 instances, as implemented in `data_loader.py`. The performance of the classifiers on this dataset is summarized in Table 4.

Table 4: Cover Type dataset: classification performance summary.

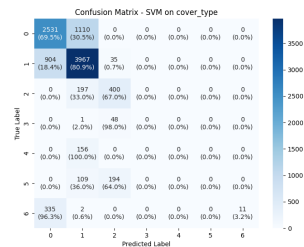
Classifier	CV Accuracy	Accuracy (full)
Naive Bayes	0.4579 ± 0.0007	0.46
KNN	0.9595 ± 0.0006	0.97
SVM	0.6957 ± 0.0041	0.69
MLP	0.7625 ± 0.0049	0.76
Random Forest	0.9466 ± 0.0012	0.95
XGBoost	0.8386 ± 0.0038	0.84



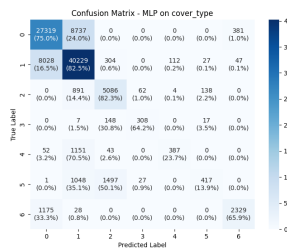
(a) Naive Bayes



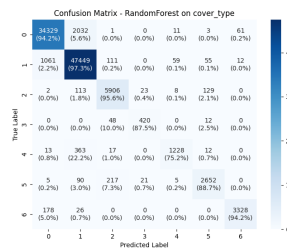
(b) KNN



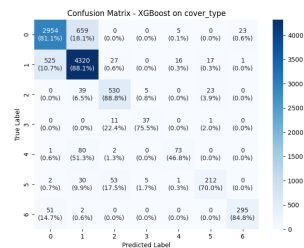
(c) SVM



(d) MLP



(e) Random Forest



(f) XGBoost

Figure 6: Confusion matrices for the cover-type dataset.

3 Analysis

After experimenting with six classifiers on four datasets, we observed several consistent patterns, as well as some interesting exceptions driven by specific data characteristics.

- **Overall, Ensemble Models (Random Forest and XGBoost) demonstrated the most robust and stable performance:** On the Telescope, Bankruptcy, and Dry Bean datasets, Random Forest and XGBoost significantly outperformed other models across nearly all metrics. This aligns with expectations, as these methods, by combining the predictions of many decision trees, can effectively handle high-dimensional features, reduce variance, and are insensitive to feature scaling. The low standard deviation on the Telescope dataset (std dev ≈ 0.003) and the high AUC on the Bankruptcy dataset (> 0.94) are clear testaments to their stability and accuracy.
- **The Bankruptcy dataset revealed the misleading nature of "Accuracy" on highly imbalanced data:** This was the most significant and unexpected finding. On the Bankruptcy dataset (Table 2), KNN and SVM achieved extremely high CV accuracy ($\approx 97\%$), yet their AUC scores were merely ≈ 0.55 (close to random guessing). A review of their confusion matrices (Figure 3) confirms that these models were predicting the majority class (class 0) almost exclusively. While this inflates the accuracy score, the models completely failed to identify the critical minority class (class 1). In contrast, Random Forest and XGBoost achieved the same 97
- **The performance of Naive Bayes was as expected, correlating strongly with its feature independence assumption:** Naive Bayes performed worst on datasets with a high number of features (Bankruptcy: 95 features, Covertypes: 54 features) where feature correlations likely exist. This is expected, as its fundamental assumption of "conditional independence" is likely violated. However, on the Dry Bean dataset, which has fewer features (16), Naive Bayes (76.4% CV Acc) surprisingly outperformed KNN, SVM, and MLP, demonstrating its efficiency when its underlying assumptions are reasonably met.
- **The Covertypes dataset was an exception, with KNN performing unexpectedly well:** On the large Covertypes dataset (Table 4), KNN unexpectedly emerged as the top performer with $\approx 96\%$ CV accuracy, even surpassing Random Forest ($\approx 94.7\%$). This was surprising, as KNN is often susceptible to the "curse of dimensionality" with 54 features.

A possible explanation is that despite the high dimensionality, the large number of samples (500,000) makes the feature space "dense" enough for instance-based logic to be highly effective, where a sample's nearest neighbors almost always belong to the same class.

- **MLP and SVM demonstrated sensitivity to data and hyperparameters:** The MLP was extremely unstable on the Dry Bean dataset (CV std dev 0.12) with very low final accuracy (28%). This is expected, as neural network models are highly sensitive to feature scaling, initialization, and hyperparameters (e.g., hidden layer size, learning rate), and can perform poorly without careful tuning. SVM struggled with the imbalanced Bankruptcy (AUC 0.54) and the large, multi-class Covertypes datasets, suggesting its default parameters (likely an RBF kernel) were not optimal for these complex tasks.

4 Experiments: Extracting Features

4.1 Motivation

As noted in the main analysis (Section 2.3), the Naive Bayes classifier experienced a catastrophic failure on the 95-feature Bankruptcy dataset. Its performance (CV ACC 0.06) was far worse than all other models. We hypothesized this was not a failure of the model itself, but a failure caused by the high-dimensional, redundant features violating the model's core assumption of feature independence.

4.2 Method

We used Random Forest, one of the top-performing classifiers (AUC 0.9448), as a feature selection tool. We first trained a Random Forest model on the full 95-feature dataset and then extracted the `feature_importances_` attribute. We selected the indices of the top 20 most important features and created a new, reduced dataset ('bankruptcy_top20'). Finally, we re-trained and re-evaluated the Naive Bayes classifier using only this 20-feature subset.

4.3 Results and Analysis

The results were dramatic and confirmed our hypothesis. When trained on the reduced 20-feature subset, the Naive Bayes classifier's performance surged across all key metrics.

Table 5 provides a direct comparison of the model’s performance before and after feature selection. The Area Under Curve (AUC), the most reliable metric for this imbalanced dataset, improved from a near-random 0.6521 to a strong 0.9027.

Table 5: Performance comparison of Naive Bayes on the Bankruptcy dataset.

Naive Bayes Model	CV Accuracy	Accuracy (full)	AUC
Full 95 Features	0.0631 ± 0.0069	0.06	0.6521
Top 20 (from RF)	0.7846 ± 0.3655	0.97	0.9027

This experiment strongly indicates that the model itself is effective, but it was crippled by the 75 redundant or noisy features in the original dataset. This highlights the sensitivity of Naive Bayes to its assumptions and demonstrates the utility of ensemble methods as powerful tools for feature selection.

5 Conclusion

This project’s experiments revealed two critical insights beyond simple classifier benchmarking. First, the Bankruptcy dataset experiments proved that for imbalanced data, accuracy is a misleading metric, and AUC is a far more reliable measure of a model’s true predictive skill. Second, our feature selection experiment confirmed that a model’s performance is fundamentally tied to its underlying assumptions. By reducing 95 redundant features to the top 20, the Naive Bayes classifier’s AUC surged from 0.6521 to 0.9027, validating this principle. These results show that successful classification depends not just on algorithm choice but on the critical evaluation of metrics and feature relevance.

A Program Listing

A.1 classifier.py

```
1  import numpy as np
2  from loguru import logger
3  from typing import Tuple
4
5  from sklearn.naive_bayes import GaussianNB
6  from sklearn.neighbors import KNeighborsClassifier
7  from sklearn.ensemble import RandomForestClassifier
8  from sklearn.neural_network import MLPClassifier
9  from sklearn.svm import SVC
10 from xgboost import XGBClassifier
11
12 class MLP:
13     def __init__(self):
14         self.model = MLPClassifier()
15         logger.info("Initialized Multi-layer Perceptron Classifier.")
16
17     def train(self, X: np.ndarray, y: np.ndarray) -> None:
18         self.model.fit(X, y)
19         logger.info("Model training completed.")
20
21     def predict(self, X: np.ndarray) -> Tuple[np.ndarray, np.ndarray]:
22         y_pred = self.model.predict(X)
23         y_scores = self.model.predict_proba(X)
24         logger.info("Prediction completed.")
25         return y_pred, y_scores
26
27 class NaiveBayes:
28     def __init__(self):
29         self.model = GaussianNB()
30         logger.info("Initialized Gaussian Naive Bayes Classifier.")
31
32     def train(self, X: np.ndarray, y: np.ndarray) -> None:
33         self.model.fit(X, y)
34         logger.info("Model training completed.")
35
```

```
36     def predict(self, X: np.ndarray) -> Tuple[np.ndarray, np.ndarray]:
37         y_pred = self.model.predict(X)
38         y_scores = self.model.predict_proba(X)
39         logger.info("Prediction completed.")
40         return y_pred, y_scores
41
42     class KNN:
43         def __init__(self):
44             self.model = KNeighborsClassifier()
45             logger.info("Initialized K-Nearest Neighbors Classifier.")
46
47         def train(self, X: np.ndarray, y: np.ndarray) -> None:
48             self.model.fit(X, y)
49             logger.info("Model training completed.")
50
51         def predict(self, X: np.ndarray) -> Tuple[np.ndarray, np.ndarray]:
52             y_pred = self.model.predict(X)
53             y_scores = self.model.predict_proba(X)
54             logger.info("Prediction completed.")
55             return y_pred, y_scores
56
57
58     class RandomForest:
59         def __init__(self):
60             self.model = RandomForestClassifier()
61             logger.info("Initialized Random Forest Classifier.")
62
63         def train(self, X: np.ndarray, y: np.ndarray) -> None:
64             self.model.fit(X, y)
65             logger.info("Model training completed.")
66
67         def predict(self, X: np.ndarray) -> Tuple[np.ndarray, np.ndarray]:
68             y_pred = self.model.predict(X)
69             y_scores = self.model.predict_proba(X)
70             logger.info("Prediction completed.")
71             return y_pred, y_scores
72
73
74     class SVM:
```

```
75     def __init__(self):
76         self.model = SVC(probability=True)
77         logger.info("Initialized Support Vector Machine Classifier.")
78
79     def train(self, X: np.ndarray, y: np.ndarray) -> None:
80         self.model.fit(X, y)
81         logger.info("Model training completed.")
82
83     def predict(self, X: np.ndarray) -> Tuple[np.ndarray, np.ndarray]:
84         y_pred = self.model.predict(X)
85         y_scores = self.model.predict_proba(X)
86         logger.info("Prediction completed.")
87         return y_pred, y_scores
88
89
90 class XGBoost:
91     def __init__(self):
92         self.model = XGBClassifier(use_label_encoder=False, eval_metric='logloss')
93         logger.info("Initialized XGBoost Classifier.")
94
95     def train(self, X: np.ndarray, y: np.ndarray) -> None:
96         self.model.fit(X, y)
97         logger.info("Model training completed.")
98
99     def predict(self, X: np.ndarray) -> Tuple[np.ndarray, np.ndarray]:
100         y_pred = self.model.predict(X)
101         y_scores = self.model.predict_proba(X)
102         logger.info("Prediction completed.")
103         return y_pred, y_scores
```

A.2 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(X: pd.DataFrame, y: pd.DataFrame) -> Tuple[np.ndarray, np.ndarray]:
45     X = X.values
```

```
46     y = y.values
47
48     le = LabelEncoder()
49     y = le.fit_transform(y.ravel())
50
51     return X, y
52
53
54 @memory.cache
55 def load_telescope() -> pd.DataFrame:
56     logger.info("Loading Telescope (id=159) from ucimlrepo...")
57
58     telescope = fetch_ucirepo(id=159)
59
60     X = telescope.data.features
61     y = telescope.data.targets
62
63     return _preprocess_data(X, y)
64
65
66 @memory.cache
67 def load_bankruptcy() -> pd.DataFrame:
68     logger.info("Loading Bank Bankruptcy (id=572) from ucimlrepo...")
69
70     bank_bankruptcy = fetch_ucirepo(id=572)
71
72     X = bank_bankruptcy.data.features
73     y = bank_bankruptcy.data.targets
74
75     return _preprocess_data(X, y)
76
77
78 @memory.cache
79 def load_cover_type() -> pd.DataFrame:
80     logger.info("Loading Cover Type (id=31) from ucimlrepo...")
81
82     cover_type = fetch_ucirepo(id=31)
83
84     X = cover_type.data.features
```

```

85     y = cover_type.data.targets
86
87     sample = np.random.choice(len(X), size=int(5E5), replace=False)
88     X = X.iloc[sample]
89     y = y.iloc[sample]
90
91     return _preprocess_data(X, y)
92
93
94 @memory.cache
95 def load_dry_bean() -> pd.DataFrame:
96     logger.info("Loading Dry Bean Classification (id=602) from ucimlrepo...")
97
98     dry_bean = fetch_ucirepo(id=602)
99
100    X = dry_bean.data.features
101    y = dry_bean.data.targets
102
103    return _preprocess_data(X, y)

```

A.3 evaluation.py

```

1  import numpy as np
2  from sklearn.metrics import confusion_matrix, classification_report, roc_curve, auc
3  from loguru import logger
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6  import os
7  from typing import Optional
8
9  def _plot_confusion_matrix(cm: np.ndarray,
10                           dataset_name: str,
11                           classifier_name: str,
12                           results_dir: str):
13
14     """
15     Plot and save confusion matrix heatmap.
16     """
17     try:

```



```
17     plt.figure(figsize=(8, 6))
18
19     num_classes = cm.shape[0]
20
21     show_percent = True
22     if num_classes > 10:
23         show_percent = False
24
25     cm_percent = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
26
27     annot_labels = (
28         np.asarray([
29             f"{val}\n({perc:.1%})" if show_percent else f"{val}"
30             for val, perc in zip(cm.flatten(), cm_percent.flatten())
31         ])
32     ).reshape(cm.shape)
33
34     sns.heatmap(
35         cm,
36         annot=annot_labels,
37         fmt='',
38         cmap='Blues',
39         cbar=True
40     )
41
42     plt.xlabel('Predicted Label')
43     plt.ylabel('True Label')
44     plt.title(f'Confusion Matrix - {classifier_name} on {dataset_name}')
45
46     os.makedirs(results_dir, exist_ok=True)
47     save_path = os.path.join(results_dir, f"CM_{dataset_name}_{classifier_name}.png")
48     plt.savefig(save_path)
49     logger.info(f"Confusion Matrix heatmap saved to {save_path}")
50     plt.close()
51
52 except Exception as e:
53     logger.error("Failed to plot confusion matrix heatmap.")
54     logger.exception(e)
55
```

```
56
57 def evaluate_classifier(y_true: np.ndarray,
58                         y_pred: np.ndarray,
59                         y_scores: np.ndarray,
60                         dataset_name: str,
61                         classifier_name: str,
62                         results_dir: str = "results") -> None:
63     """
64     Evaluate classifier performance and generate reports.
65
66     Args:
67         y_true (np.ndarray): Ground truth labels
68         y_pred (np.ndarray): Predicted labels from the model
69         y_scores (np.ndarray): Model discriminant function values (e.g., probabilities)
70         dataset_name (str): Dataset name (used for logs and filenames)
71         classifier_name (str): Classifier name (used for logs and filenames)
72         results_dir (str): Directory to save plots (e.g., ROC curve)
73     """
74
75     logger.info(f"--- Evaluation for [{classifier_name}] on [{dataset_name}] ---")
76
77     try:
78         logger.info("Confusion Matrix:")
79         cm = confusion_matrix(y_true, y_pred)
80
81         _plot_confusion_matrix(cm, dataset_name, classifier_name, results_dir)
82
83         logger.info(f"\n{cm}")
84
85         logger.info("Classification Report:")
86         report = classification_report(y_true, y_pred, zero_division=0)
87         logger.info(f"\n{report}")
88
89         num_classes = len(np.unique(y_true))
90
91         if num_classes == 2:
92             logger.info("Two-class dataset detected. Calculating ROC/AUC...")
93
94             scores_for_roc = y_scores[:, 1]
```

```

95
96         fpr, tpr, thresholds = roc_curve(y_true, scores_for_roc)
97
98         roc_auc = auc(fpr, tpr)
99         logger.success(f"Area Under Curve (AUC): {roc_auc:.4f}")
100
101         plt.figure(figsize=(8, 6))
102         plt.plot(fpr, tpr, color='darkorange', lw=2,
103                 label=f'ROC curve (area = {roc_auc:.4f})')
104         plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
105         plt.xlim([0.0, 1.0])
106         plt.ylim([0.0, 1.05])
107         plt.xlabel('False Positive Rate')
108         plt.ylabel('True Positive Rate')
109         plt.title(f'ROC Curve - {classifier_name} on {dataset_name}')
110         plt.legend(loc="lower right")
111
112         os.makedirs(results_dir, exist_ok=True)
113         save_path = os.path.join(results_dir, f"ROC_{dataset_name}_{classifier_name}.png")
114         plt.savefig(save_path)
115         logger.info(f"ROC curve saved to {save_path}")
116         plt.close()
117
118     else:
119         logger.info(f"Skipping ROC/AUC calculation ({num_classes} classes detected).")
120
121     except Exception as e:
122         logger.error("An error occurred during evaluation.")
123         logger.exception(e)
124         raise
125
126     logger.info(f"--- End of Evaluation for [{classifier_name}] ---")

```

A.4 main.py

```

1  import argparse
2  import numpy as np
3  from loguru import logger

```

```
4  from sklearn.model_selection import KFold
5  from sklearn.metrics import accuracy_score
6  import sys
7
8  from data_loader import load_dataset
9  from evaluation import evaluate_classifier
10
11 from classifier import MLP, NaiveBayes, KNN, RandomForest, SVM, XGBoost
12
13 CLASSIFIER_MAP = {
14     'MLP': MLP,
15     'NaiveBayes': NaiveBayes,
16     'KNN': KNN,
17     'RandomForest': RandomForest,
18     'SVM': SVM,
19     'XGBoost': XGBoost
20 }
21
22
23 def k_fold_cross_validation(classifier_class, X, y, k=5, random_state=42):
24     """
25     Perform k-fold cross validation on the given classifier.
26
27     Args:
28         classifier_class (class): The classifier class to instantiate.
29         X (np.ndarray): Feature matrix
30         y (np.ndarray): Target labels
31         k (int): Number of folds
32         random_state (int): Random seed for reproducibility
33
34     Returns:
35         dict: Dictionary containing fold results and average metrics
36     """
37     logger.info(f"Starting {k}-fold cross validation for {classifier_class.__name__}...")
38
39     kfold = KFold(n_splits=k, shuffle=True, random_state=random_state)
40
41     fold_results = []
42
```

```
43     for fold_idx, (train_idx, val_idx) in enumerate(kfold.split(X), 1):
44         logger.info(f"\n{'='*60}")
45         logger.info(f"Fold {fold_idx}/{k}")
46         logger.info(f"{'='*60}")
47
48         # Split data for this fold
49         X_train, X_val = X[train_idx], X[val_idx]
50         y_train, y_val = y[train_idx], y[val_idx]
51
52         logger.info(
53             f"Train size: {len(X_train)}, Validation size: {len(X_val)}")
54
55         classifier = classifier_class()
56         classifier.train(X_train, y_train)
57
58         # Make predictions
59         y_pred, y_scores = classifier.predict(X_val)
60
61         # Calculate accuracy
62         accuracy = accuracy_score(y_val, y_pred)
63         logger.success(f"Fold {fold_idx} Accuracy: {accuracy:.4f}")
64
65         # Store results
66         fold_results.append({
67             'fold': fold_idx,
68             'accuracy': accuracy,
69             'y_true': y_val,
70             'y_pred': y_pred,
71             'y_scores': y_scores
72         })
73
74         # Calculate average metrics
75         accuracies = [result['accuracy'] for result in fold_results]
76         avg_accuracy = np.mean(accuracies)
77         std_accuracy = np.std(accuracies)
78
79         logger.info(f"\n{'='*60}")
80         logger.info("Cross Validation Results Summary")
81         logger.info(f"{'='*60}")
```

```
82     logger.success(f"Average Accuracy: {avg_accuracy:.4f} (+/- {std_accuracy:.4f})")
83     logger.info(f"Individual fold accuracies: {[f'{acc:.4f}' for acc in accuracies]}")
84
85     return {
86         'fold_results': fold_results,
87         'avg_accuracy': avg_accuracy,
88         'std_accuracy': std_accuracy,
89         'accuracies': accuracies
90     }
91
92
93 if __name__ == "__main__":
94     parser = argparse.ArgumentParser(description='Test Classifier with k-fold cross validation')
95     parser.add_argument('--model', type=str, required=True,
96                         help=f'Classifier to use. Available: {list(CLASSIFIER_MAP.keys())}')
97     parser.add_argument('--dataset', type=str, required=True,
98                         help='Dataset to load (e.g., iris, breast_cancer, bank_note, glass)')
99     parser.add_argument('--k', type=int, default=5,
100                        help='Number of folds for cross validation (default: 5)')
101     parser.add_argument('--random_state', type=int, default=42,
102                        help='Random seed for reproducibility (default: 42)')
103     args = parser.parse_args()
104
105     # Load dataset
106     X, y = load_dataset(args.dataset)
107     logger.info(f"Loaded {args.dataset} dataset: X shape={X.shape}, y shape={y.shape}")
108
109     classifier_class = CLASSIFIER_MAP.get(args.model)
110     if classifier_class is None:
111         logger.error(f"Unknown model: '{args.model}'")
112         logger.info(f"Available models: {list(CLASSIFIER_MAP.keys())}")
113         sys.exit(1)
114
115     logger.info(f"Selected model: {args.model}")
116
117     cv_results = k_fold_cross_validation(
118         classifier_class, X, y, k=args.k, random_state=args.random_state
119     )
120
```

```
121     logger.info(f"\n{'='*60}")
122     logger.info("Training final model on entire dataset for evaluation...")
123     logger.info(f"{'='*60}")
124
125     classifier = classifier_class()
126     classifier.train(X, y)
127     y_pred_final, y_scores_final = classifier.predict(X)
128
129     evaluate_classifier(
130         y_true=y,
131         y_pred=y_pred_final,
132         y_scores=y_scores_final,
133         dataset_name=args.dataset,
134         classifier_name=args.model
135     )
```

A.5 subset.py

```
1  from sklearn.ensemble import RandomForestClassifier
2  import numpy as np
3  from loguru import logger
4
5  from data_loader import load_dataset
6  from classifier import NaiveBayes
7  from main import k_fold_cross_validation, evaluate_classifier
8
9
10 if __name__ == "__main__":
11     dataset = "bankruptcy"
12
13     model = RandomForestClassifier()
14     logger.info("Initialized Random Forest Classifier.")
15
16     X, y = load_dataset(dataset)
17     model.fit(X, y)
18
19     importances = model.feature_importances_
20     top_20_indices = np.argsort(importances)[::-1][:20]
```

```
21
22     logger.info(f"Top 20 most important feature indices: {top_20_indices}")
23
24     X_top20 = X[:, top_20_indices]
25     logger.info(f"Created new dataset with top 20 features: X shape={X_top20.shape}")
26
27     logger.info("Running K-fold CV for NaiveBayes on TOP 20 FEATURES...")
28     nb_top20_cv_results = k_fold_cross_validation(
29         NaiveBayes, X_top20, y
30     )
31
32     logger.info("Evaluating final NaiveBayes model on TOP 20 FEATURES...")
33     nb_top20_final = NaiveBayes()
34     nb_top20_final.train(X_top20, y)
35     y_pred_nb_top20, y_scores_nb_top20 = nb_top20_final.predict(X_top20)
36
37     evaluate_classifier(
38         y_true=y,
39         y_pred=y_pred_nb_top20,
40         y_scores=y_scores_nb_top20,
41         dataset_name="bankruptcy_top20",
42         classifier_name="NaiveBayes"
43     )
```
