

```

1  import numpy as np
2  import os
3  from sklearn.model_selection import GridSearchCV
4  from xgboost import XGBClassifier
5  from sklearn.metrics import accuracy_score, classification_report
6  import json
7
8  # 전처리한 파일 이용
9  base_path = r"C:\Users\minar\OneDrive\바탕 화면\학교\기계학습\팀프로젝트\dataset"
10 X_train_path = os.path.join(base_path, "X_train.npy")
11 y_train_path = os.path.join(base_path, "y_train.npy")
12 X_test_path = os.path.join(base_path, "X_test.npy")
13 y_test_path = os.path.join(base_path, "y_test.npy")
14
15 # Numpy 배열 로드
16 X_train_np = np.load(X_train_path)
17 y_train_np = np.load(y_train_path)
18 X_test_np = np.load(X_test_path)
19 y_test_np = np.load(y_test_path)
20
21 # XGBoost 모델 정의
22 model = XGBClassifier(
23     objective='multi:softmax',
24     num_class=11, #전체 클래스 갯수 11개
25     eval_metric='mlogloss',
26     verbosity=1
27 )

```

1. 데이터 로드

- 전처리된 학습 데이터와 테스트 데이터를 Numpy 형식으로 불러옵니다.
- 입력 데이터 X_train, X_test와 레이블 y_train, y_test를 각각 로드합니다.

2. 모델 초기화

- **XGBoost** 모델을 초기화하며, 다중 클래스 분류를 위해 ****objective='multi:softmax'****를 설정했습니다.
- **num_class=11**: 클래스 개수가 11개임을 명시하였습니다.
- **eval_metric='mlogloss'**: 다중 클래스 분류에서 손실 함수로 **로그 손실(Log Loss)**을 사용합니다.

```

1  # 하이퍼파라미터 그리드 정의
2  param_grid = {
3      'max_depth': [3, 6],
4      'eta': [0.1, 0.05, 0.01],
5      'n_estimators': [50, 100, 150],
6  }
7
8  # GridSearchCV 정의
9  grid_search = GridSearchCV(
10     estimator=model,
11     param_grid=param_grid,
12     scoring='accuracy',
13     cv=3,
14     verbose=1,
15     n_jobs=-1,
16     return_train_score=True
17 )
18
19 # 그리드 서치 학습
20 print("GridSearchCV 시작...")
21 grid_search.fit(X_train_np, y_train_np)
22
23 # 최적 하이퍼파라미터 출력
24 print(f"최적 하이퍼파라미터: {grid_search.best_params_}")
25 print(f"Best Cross-Validated Accuracy: {grid_search.best_score_:.2f}")
26

```

1. 하이퍼파라미터 탐색

- XGBoost의 주요 하이퍼파라미터를 그리드 탐색(Grid Search)을 통해 최적화합니다.
- 탐색할 파라미터:
 - max_depth: 트리의 최대 깊이 (3, 6)
 - eta: 학습률 (0.1, 0.05, 0.01)
 - n_estimators: 트리 개수 (50, 100, 150)

2. GridSearchCV

- 교차 검증을 사용해 모든 파라미터 조합을 평가합니다.
- **cv=3**: 3-fold 교차 검증을 수행합니다.

- **scoring='accuracy'**: 정확도를 기준으로 최적의 조합을 찾습니다.

결과 출력

- 최적의 하이퍼파라미터와 ****검증 정확도(Mean CV Accuracy)****를 출력합니다.

```

1 def print_detailed_cv_results(grid_search):
2     cv_results = grid_search.cv_results_
3     print("\nDetailed Results for Each Hyperparameter Combination:")
4     for i in range(len(cv_results['params'])):
5         print(f"Hyperparameters: {cv_results['params'][i]}")
6         print(f"Mean Train Accuracy: {cv_results['mean_train_score'][i]:.2f}")
7         print(f"Mean Validation Accuracy: {cv_results['mean_test_score'][i]:.2f}")
8         for fold in range(grid_search.cv):
9             train_score = cv_results[f'split{fold}_train_score'][i]
10            test_score = cv_results[f'split{fold}_test_score'][i]
11            print(f"  Fold {fold}: Train Accuracy = {train_score:.2f}, Validation Accuracy = {test_score:.2f}")
12        print("-" * 50)
13
14 print_detailed_cv_results(grid_search)
15
16 # 최적 모델로 테스트 세트 예측
17 best_model = grid_search.best_estimator_
18 y_pred = best_model.predict(X_test_np)
19
20 # 학습 데이터 정확도
21 y_train_pred = best_model.predict(X_train_np)
22 train_accuracy = accuracy_score(y_train_np, y_train_pred)
23 print(f"Train Accuracy: {train_accuracy * 100:.2f}%")
24
25 # 정확도 및 분류 보고서
26 test_accuracy = accuracy_score(y_test_np, y_pred)
27 print(f"Test Accuracy: {test_accuracy * 100:.2f}%")

```

1. 세부 결과 출력

- `cv_results_`를 이용해 각 하이퍼파라미터 조합의 **Train/Validation Accuracy**를 출력합니다.
- 각 Fold별 교차 검증 결과도 함께 출력하여 모델 성능을 세부적으로 확인합니다.

2. 최적 모델 예측 및 평가

- **최적의 하이퍼파라미터**를 적용한 모델로 테스트 데이터 예측을 수행합니다.
- 학습 정확도(Train Accuracy)와 테스트 정확도(Test Accuracy)를 계산하여 출력합니다.

```

1 save_path = os.path.join(base_path, "xgboost_best_model.json")
2 best_model.save_model(save_path)
3 print(f"Best model saved to {save_path}")
4
5 # 메타데이터 생성, 실험 결과에 대한 간략한 데이터
6 cv_results = grid_search.cv_results_
7 model_metadata = {
8     "test_accuracy": test_accuracy,
9     "train_accuracy": train_accuracy,
10    "classification_report": classification_report(y_test_np, y_pred, output_dict=True),
11    "best_hyperparameters": grid_search.best_params_,
12    "best_cv_accuracy": grid_search.best_score_,
13    "cv_results": [
14        {
15            "params": cv_results['params'][i],
16            "mean_train_score": cv_results['mean_train_score'][i],
17            "mean_test_score": cv_results['mean_test_score'][i],
18            "fold_scores": {
19                f"fold_{fold}": {
20                    "train": cv_results[f'split{fold}_train_score'][i],
21                    "validation": cv_results[f'split{fold}_test_score'][i]
22                }
23                for fold in range(grid_search.cv)
24            }
25        }
26        for i in range(len(cv_results['params']))
27    ]
28 }
29
30 # 메타데이터 저장
31 metadata_save_path = os.path.join(base_path, "xgboost_best_model_metadata.json")
32 with open(metadata_save_path, "w") as f:
33     json.dump(model_metadata, f, indent=4)
34 print(f"Metadata saved to {metadata_save_path}")
35

```

1. 모델 저장

- 최적의 하이퍼파라미터로 학습된 모델을 **JSON 형식**으로 저장합니다.
- 저장된 모델은 나중에 불러와 바로 예측에 사용할 수 있습니다.

2. 메타데이터 저장

- 실험 결과를 기록한 **메타데이터 파일**을 JSON 형식으로 저장합니다.
- 기록된 내용:
 - 테스트 및 학습 정확도
 - 최적 하이퍼파라미터
 - 교차 검증 결과

- 분류 보고서

3. 결과 확인

- 저장된 모델 및 메타데이터 파일은 다른 실험과 비교하거나 후속 분석에 활용할 수 있습니다.