

모형진단과 교차검증 (실습)

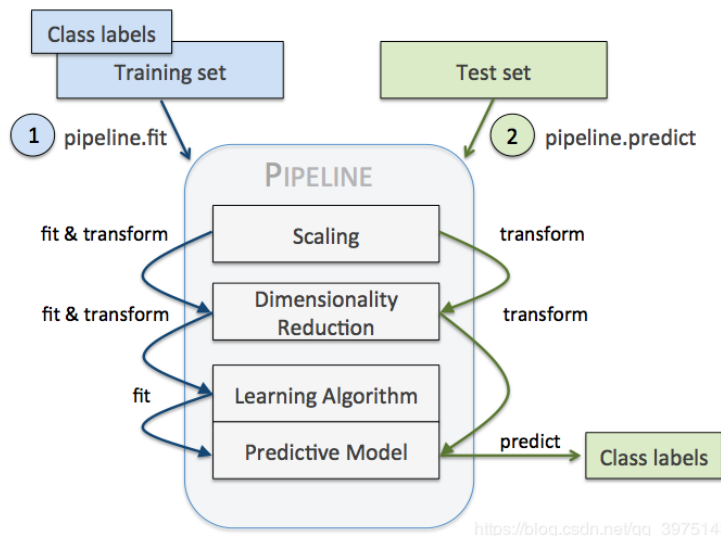
Evaluation & Cross Validation

파이프라인 실습 #1

■ 파이프 라인(Pipeline)

- 사이킷런의 Pipeline 클래스는 연속된 변환을 순차적으로 처리할 수 있는 기능을 제공하는 유용한 래퍼(Wrapper) 도구

```
pipe_lr = make_pipeline(StandardScaler(),  
                        PCA(n_components=2),  
                        LogisticRegression(solver='liblinear', random_state=1))  
  
pipe_lr.fit(X_train, y_train)  
y_pred = pipe_lr.predict(X_test)  
print('테스트 정확도: %.3f' % pipe_lr.score(X_test, y_test))
```



파이프라인 실습 #1

- 데이터셋: 유방암 데이터
- 학습/시험 데이터: X, 학습/시험 데이터 라벨: Y

① 데이터 로드

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load\_breast\_cancer.html
2 # Dimensionality: 30, Classes: 2
3 # 212(M-유방암)-label('0'), 357(B-정상인) - label('1')
4 from sklearn.datasets import load_breast_cancer
5 cancer = load_breast_cancer()
6 X = cancer.data
7 Y = cancer.target
```

② 데이터 분할 (학습/테스트)

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.model\_selection.train\_test\_split.htm
2
3 # 학습 데이터 분할
4 from sklearn.model_selection import train_test_split
5 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, stratify=Y, random_state=1)
```

파이프라인 실습 #1

- 데이터셋: 유방암 데이터
- 학습/시험 데이터: X, 학습/시험 데이터 라벨: Y

③ 파이프라인 모듈 설계

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.pipeline.Pipeline.html
2 # 파이프라인 기능을 이용한 모듈 설계
3
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.decomposition import PCA
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.pipeline import make_pipeline
8
9 pipeline = make_pipeline(StandardScaler(), PCA(n_components=4), LogisticRegression() )
```

④ 모델 학습 및 평가

```
1 # 모델 학습 및 테스트
2 pipeline.fit(X_train, Y_train)
3 Y_Pred = pipeline.predict(X_test)
```

0.9736842105263158

```
1 # 모델 평가
2 from sklearn.metrics import accuracy_score
3 accuracy_score(Y_test, Y_Pred)
```

0.9736842105263158

교차검증 실습 #1

- 데이터셋: 유방암 데이터
- 학습/시험 데이터: X, 학습/시험 데이터 라벨: Y

① 데이터 로드

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load\_breast\_cancer.html
2 # Dimensionality: 30, Classes: 2
3 # 212(M-유방암)-label('0'), 357(B-정상인) - label('1')
4 from sklearn.datasets import load_breast_cancer
5 cancer = load_breast_cancer()
6 X = cancer.data
7 Y = cancer.target
```

② 데이터 분할 (학습/테스트)

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.model\_selection.train\_test\_split.htm
2
3 # 학습 데이터 분할
4 from sklearn.model_selection import train_test_split
5 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, stratify=Y, random_state=1)
```

교차검증 실습 #1

- 데이터셋: 유방암 데이터
- 학습/시험 데이터: X, 학습/시험 데이터 라벨: Y

③ 파이프라인 모듈 설계

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.pipeline.Pipeline.html
2 # 파이파라인 기능을 이용한 모듈 설계
3
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.decomposition import PCA
6 from sklearn.linear_model import LogisticRegression
7 from sklearn.pipeline import make_pipeline
8
9 pipeline = make_pipeline(StandardScaler(), PCA(n_components=4), LogisticRegression() )
```

④ 교차검증

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.model\_selection.cross\_validate.html
2
3 from sklearn.model_selection import cross_validate
4 scores = cross_validate(pipeline, X_train, Y_train, cv=10, return_train_score=True)
```

교차검증 실습 #1

- 데이터셋: 유방암 데이터
- 학습/시험 데이터: X, 학습/시험 데이터 라벨: Y

⑤ 과적합 여부 분석

```
1 import numpy as np
2
3 print('CV Validation Accuracy scores: ', scores['train_score'])
4 print('CV Validation Accuracy: %.3f +/- %.3f' %(np.mean(scores['train_score']), np.std(scores['train_score'])))
```

```
CV Validation Accuracy scores: [0.96577017 0.96577017 0.96577017 0.96577017 0.96821516 0.96585366
0.97073171 0.96829268 0.97560976 0.96585366]
CV Validation Accuracy: 0.968 +/- 0.003
```

```
1 import numpy as np
2
3 print('CV Validation Accuracy scores: ', scores['test_score'])
4 print('CV Validation Accuracy: %.3f +/- %.3f' %(np.mean(scores['test_score']), np.std(scores['test_score'])))
```

```
CV Validation Accuracy scores: [0.97826087 0.97826087 0.95652174 1.          0.95652174 0.97777778
0.93333333 0.95555556 0.91111111 1.          ]
CV Validation Accuracy: 0.965 +/- 0.027
```

교차검증 실습 #1

- 데이터셋: 유방암 데이터
- 학습/시험 데이터: X, 학습/시험 데이터 라벨: Y

⑥ 최적 모델 탐색

```
1 from sklearn.model_selection import GridSearchCV
2
3 parameters = {}
4 gs = GridSearchCV(pipeline, parameters, scoring='accuracy', cv=10)
5 gs.fit(X_train, Y_train)
```

```
1 best = gs.best_estimator_
```

```
1 gs.cv_results_
```

```
{'mean_fit_time': array([0.00662093]),
 'mean_score_time': array([0.00051041]),
 'mean_test_score': array([0.9647343]),
 'params': [{}],
 'rank_test_score': array([1], dtype=int32),
 'split0_test_score': array([0.97826087]),
 'split1_test_score': array([0.97826087]),
 'split2_test_score': array([0.95652174]),
 'split3_test_score': array([1.]),
 'split4_test_score': array([0.95652174]),
 'split5_test_score': array([0.97777778]),
 'split6_test_score': array([0.93333333]),
 'split7_test_score': array([0.95555556]),
 'split8_test_score': array([0.91111111]),
 'split9_test_score': array([1.]),
 'std_fit_time': array([0.00164769]),
 'std_score_time': array([3.95498696e-05]),
 'std_test_score': array([0.02665336])}
```


교차검증 실습 #1

- 데이터셋: 유방암 데이터
- 학습/시험 데이터: X , 학습/시험 데이터 라벨: Y

⑦ 최적 모델 평가

```
1 from sklearn.metrics import accuracy_score
2
3 Y_train_Pred = best.predict(X_train)
4 accuracy_score(Y_train, Y_train_Pred)
5
```

0.967032967032967

```
1 from sklearn.metrics import accuracy_score
2
3 Y_test_Pred = best.predict(X_test)
4 accuracy_score(Y_test, Y_test_Pred)
5
```

0.9736842105263158

최적 모델 찾기 실습 #1

- 데이터셋: 유방암 데이터
- 학습/시험 데이터: X, 학습/시험 데이터 라벨: Y

① 데이터 로드

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.datasets.load\_breast\_cancer.html
2 # Dimensionality: 30, Classes: 2
3 # 212(M-유방암)-label('0'), 357(B-정상인) - label('1')
4 from sklearn.datasets import load_breast_cancer
5 cancer = load_breast_cancer()
6 X = cancer.data
7 Y = cancer.target
```

② 데이터 분할 (학습/테스트)

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.model\_selection.train\_test\_split.htm
2
3 # 학습 데이터 분할
4 from sklearn.model_selection import train_test_split
5 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2, stratify=Y, random_state=1)
```

최적 모델 찾기 실습 #1

- 데이터셋: 유방암 데이터
- 학습/시험 데이터: X, 학습/시험 데이터 라벨: Y

③ 파이프라인 모듈 설계

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.pipeline.Pipeline.html
2 # 파이파라인 기능을 이용한 모듈 설계
3
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.svm import SVC
6 from sklearn.pipeline import Pipeline
7 from sklearn.model_selection import GridSearchCV
8
9 # 파라미터 Parsing
10 estimators = [('normalization', StandardScaler()), ('clf', SVC())]
11 pipe = Pipeline(estimators)
```

④ GridSearch 파라미터 세팅

```
1 # https://scikit-learn.org/stable/modules/generated/sklearn.model\_selection.ParameterGrid.html#sklearn.model\_selection.ParameterGrid
2
3 from sklearn.model_selection import ParameterGrid
4 grid = [{'clf__kernel': ['linear'], 'clf__C': [0.001, 0.01, 0.1, 1, 10, 100, 1000]},
5         {'clf__kernel': ['rbf'], 'clf__gamma': [0.001, 0.01, 0.1, 1, 10, 100, 1000], 'clf__C': [0.001, 0.01, 0.1, 1, 10, 100, 1000]}]
6
7 grid_param = ParameterGrid(grid)
8 list(grid_param)
```

최적 모델 찾기 실습 #1

- 데이터셋: 유방암 데이터
- 학습/시험 데이터: X, 학습/시험 데이터 라벨: Y

⑤ GridSearch 정의 및 최적 모델 탐색

```
1 gs = GridSearchCV(pipe, grid_param, scoring='accuracy', cv=10, n_jobs=1)
```

```
1 gs.fit(X_train, Y_train)
```

```
1 print(gs.best_score_)
```

```
0.9758454106280192
```

```
1 print(gs.best_params_)
```

```
{'clf__C': 10, 'clf__gamma': 0.001, 'clf__kernel': 'rbf'}
```

⑥ 최적 모델 평가

```
1 best_model = gs.best_estimator_  
2 Y_test_pred = best_model.predict(X_test)
```

```
1 from sklearn.metrics import accuracy_score  
2  
3 Y_test_Pred = best_model.predict(X_test)  
4 accuracy_score(Y_test, Y_test_Pred)
```

```
0.9824561403508771
```