

5. Scikit-Learn (sklearn)

이번 시간에는 Scikit-Learn (sklearn) 패키지를 사용하여 데이터 학습과 모델 테스트를 진행해보겠습니다.

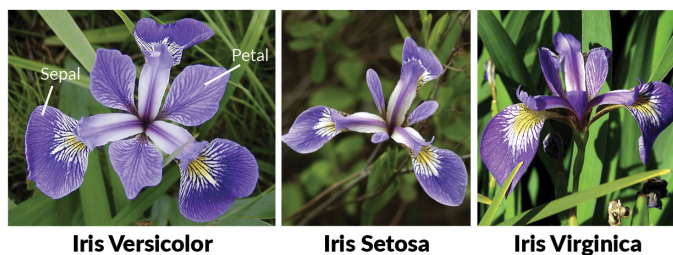
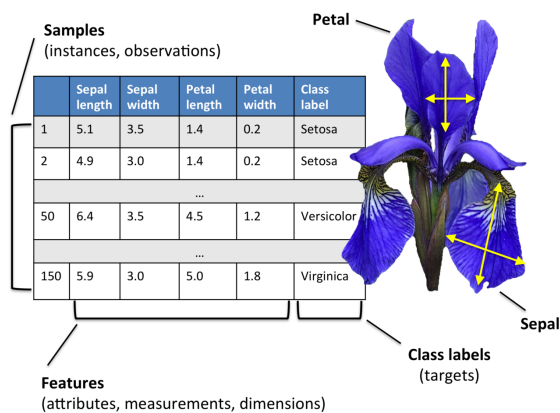
5.1. Random Forest - Ensemble of Decision Trees

먼저 필요한 라이브러리들을 불러옵니다. 오늘 사용할 라이브러리들이 다음의 code snippet에 기술되어 있습니다.

Scikit-Learn (sklearn) 패키지에는 ML 공부를 할 때 자주 사용되는 데이터셋들도 포함되어 있습니다. 그러한 데이터셋 중 Iris라는 데이터셋을 사용해보겠습니다. 이는 iris 꽃잎의 체적에 관한 정보에 근거하여 꽃의 품종을 결정하고자하는 데이터셋입니다.

Attribute Information:

1. sepal length in cm
2. sepal width in cm
3. petal length in cm
4. petal width in cm
5. class: {Setosa, Versicolour, Virginica}



```

1 # Import libraries
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.model_selection import train_test_split
4 from sklearn.model_selection import GridSearchCV
5 #from sklearn.model_selection import cross_validate
6 import warnings
7 warnings.simplefilter(action='ignore', category=FutureWarning)
8 warnings.simplefilter(action='ignore', category=DeprecationWarning)
9
10
11 # Classification metrics
12 from sklearn.metrics import (accuracy_score, precision_score,
13                             recall_score, f1_score, log_loss)
14
15 # Load the Iris dataset
16 from sklearn.datasets import load_iris
17 X, y = load_iris(return_X_y=True)
18
19 list(zip(X[:5], y[:5])) # show first 5 instances
20

```



```
[(array([5.1, 3.5, 1.4, 0.2]), 0),
 (array([4.9, 3. , 1.4, 0.2]), 0),
 (array([4.7, 3.2, 1.3, 0.2]), 0),
 (array([4.6, 3.1, 1.5, 0.2]), 0),
 (array([5. , 3.6, 1.4, 0.2]), 0)]
```

이제 수업을 통해 이야기 나눈 절차에 따라 sklearn을 사용해보겠습니다.

▼ Step 1 - Separating training and testing datasets

첫 단계로 학습된 모델의 정확한 검증을 위해 데이터셋을 training set과 testing set으로 분리를 합니다.

```
1 X_tr, X_ts, y_tr, y_ts = train_test_split(X, y, test_size=0.3, random_state=777)
```

▼ Steps 2 & 3 - Optimize hyperparameters via cross validation

먼저 시도해볼 모델은 Random Forest입니다. GridSearchCV를 사용하여 학습에 사용될 hyperparameters를 찾습니다. GridSearchCV의 fit() 함수를 사용하면, Hyperparameter 탐색과 학습 과정이 모두 수행됩니다. 다음 code snippet을 수행하고 나면, best_clf에 최종적으로 학습된 모델이 저장됩니다.

```
1 # Instantiate a model object
2 clf = RandomForestClassifier()
3
4 # Set a search range
5 parameters = {'n_estimators': [100, 150, 200],
6               'criterion': ['gini', 'entropy']}
7 # Find the best hyperparameters using GridSearchCV
8 gridsearch = GridSearchCV(clf, parameters, scoring='accuracy', cv=5)
9 gridsearch.fit(X_tr, y_tr)
10
11 # Show the best hyperparameters
12 print(f'gridsearch.best_params_ = {gridsearch.best_params_}')
13
14 # The best model is stored in 'best_clf'
15 best_clf = gridsearch.best_estimator_
16 best_clf
```

```
gridsearch.best_params_ = {'criterion': 'gini', 'n_estimators': 100}
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                        max_depth=None, max_features='auto', max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, n_estimators=100,
                        n_jobs=None, oob_score=False, random_state=None,
                        verbose=0, warm_start=False)
```

▼ Step 4 - Model performance

얻어진 모델을 앞서 떼어두었던 테스트 데이터셋에 적용해봅니다. predict() 함수를 통해 예측이 이루어지며, accuracy_score() 함수를 사용하면 쉽게 정확도를 계산해볼 수 있습니다.

```
1 y_pred = best_clf.predict(X_ts)
2 test_acc = accuracy_score(y_ts, y_pred)
3 print(f'test_acc = {test_acc}')
```

```
test_acc = 0.9777777777777777
```

▼ Step 5 - Train final model on full dataset (optional step)

선택적으로 train data와 test data를 모두 고려하여, 최종 배포용 모델을 다시 학습할 수 있습니다. 이 때, 앞서 탐색을 통해 찾은 hyperparameters를 다시 사용하면 됩니다.

```
1 final_model = RandomForestClassifier(**gridsearch.best_params_)
2 final_model.fit(X, y)
```

```
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
                        max_depth=None, max_features='auto', max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=1, min_samples_split=2,
                        min_weight_fraction_leaf=0.0, n_estimators=100,
                        n_jobs=None, oob_score=False, random_state=None,
                        verbose=0, warm_start=False)
```

▼ 5.2. Logistic Regression

이번 절에서는 Scikit-Learn (sklearn)에 포함되어있는 logistic regression 모듈을 사용하여, iris 데이터셋 학습을 진행해보겠습니다.

Logistic regression을 사용할 때, input features는 모두 같은 scale을 갖도록 해주어야 합니다. `sklearn.preprocessing.MinMaxScaler`를 사용하면, 손쉽게 normalization을 진행할 수 있습니다.

```
1 from sklearn.linear_model import LogisticRegression
2 from sklearn.preprocessing import MinMaxScaler
3
4 # (Re-)Load a dataset
5 X, y = load_iris(return_X_y=True)
6
7 # Step 1: Get training and testing datasets
8 X_tr, X_ts, y_tr, y_ts = train_test_split(X, y, test_size=0.3, random_state=777)
9
10 # Data normalization
11 normalizer = MinMaxScaler(feature_range=(0, 1))
12 normalizer.fit(X_tr)
13 X_tr_normalized = normalizer.transform(X_tr)
14 X_ts_normalized = normalizer.transform(X_ts)
15
16 # Show first 5 instances
17 print('Before normalization:\n', X_tr[:5])
18 print('After normalization:\n', X_tr_normalized[:5])
```

```
Before normalization:
[[6.7 3.3 5.7 2.5]
 [6.4 2.8 5.6 2.2]
 [5. 3.3 1.4 0.2]
 [7.7 2.6 6.9 2.3]
 [5.7 2.6 3.5 1. ]]
After normalization:
[[0.65714286 0.54166667 0.79661017 1.         ]
 [0.57142857 0.33333333 0.77966102 0.875       ]
 [0.17142857 0.54166667 0.06779661 0.04166667]
 [0.94285714 0.25         1.         0.91666667]
 [0.37142857 0.25         0.42372881 0.375       ]]
```

이제 학습과 테스트를 진행해봅니다.

```
1 # Step 2: Use GridSearchCV to find optimal hyperparameter values
2 clf = LogisticRegression()
3 parameters = {'penalty': ['l2'],
4               'C': [10e-5, 10e-3, 10e-2, 10e-1, 10e0, 10e1, 10e2, 10e3, 10e5]}
5 gridsearch = GridSearchCV(clf, parameters, scoring='accuracy', cv=5)
6 gridsearch.fit(X_tr_normalized, y_tr)
7 print(f'gridsearch.best_params_ = {gridsearch.best_params_}')
8
9 # Step 3: Get model with best hyperparameters
10 best_clf = gridsearch.best_estimator_
11
12 # Step 4: Get best model performance from testing set
13 y_pred = best_clf.predict(X_ts_normalized)
14 test_acc = accuracy_score(y_ts, y_pred)
15 print(f'test_acc = {test_acc}')
16
```

```
gridsearch.best_params_ = {'C': 100.0, 'penalty': 'l2'}
test_acc = 0.9777777777777777
```

▼ 5.3. Benchmark

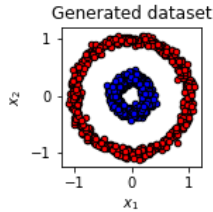
이번 절에서는 조금 복잡한 데이터셋을 생성하여 random forest와 logistic regression 모델에 각각 적용해보고, 예측 결과를 비교해보고자 합니다. 우리는 `sklearn`에 내장된 모듈을 사용하여 non-linear decision boundary를 갖는 데이터셋을 생성해보고, 이 데이터셋에 대하여 각 모델의 예측 정확도가 어떻게 되는지 비교해보고자 합니다.

▼ Dataset

`make_circles` 모듈을 사용하여 1,000개의 data instances를 생성합니다. 이어지는 코드는 `matplotlib`으로 생성된 데이터셋을 시각화하여 줍니다.

```
1 from sklearn.datasets import make_circles
2 import matplotlib.pyplot as plt
3
4 # Generate a dataset
5 X, y = make_circles(n_samples=1000, factor=.3, noise=.05)
6
7 # Show the dataset
8 plt.figure()
9 plt.subplot(2, 2, 1, aspect='equal')
10 plt.title("Generated dataset")
11 reds = y == 0
12 blues = y == 1
13
14 plt.scatter(X[reds, 0], X[reds, 1], c="red", s=20, edgecolor='k')
15 plt.scatter(X[blues, 0], X[blues, 1], c="blue", s=20, edgecolor='k')
16 plt.xlabel("$x_1$")
17 plt.ylabel("$x_2$")
```

Text(0, 0.5, '\$x_2\$')



▼ LogReg

먼저 logistic regression을 적용해봅니다. 위에서와 동일한 구조의 코드를 사용합니다.

```
1 # Step 1: Get training and testing datasets
2 X_tr, X_ts, y_tr, y_ts = train_test_split(X, y, test_size=0.3, random_state=777)
3
4 # Data normalization
5 normalizer = MinMaxScaler(feature_range=(0, 1))
6 normalizer.fit(X_tr)
7 X_tr_normalized = normalizer.transform(X_tr)
8 X_ts_normalized = normalizer.transform(X_ts)
9
10
11 # Step 2: Use GridSearchCV to find optimal hyperparameter values
12 clf = LogisticRegression()
13 parameters = {'penalty': ['l2'],
14               'C': [10e-5, 10e-3, 10e-2, 10e-1, 10e0, 10e1, 10e2, 10e3, 10e5]}
15 gridsearch = GridSearchCV(clf, parameters, scoring='accuracy', cv=5)
16 gridsearch.fit(X_tr_normalized, y_tr)
17 print(f'gridsearch.best_params_ = {gridsearch.best_params_}')
18
19 # Step 3: Get model with best hyperparameters
20 best_clf = gridsearch.best_estimator_
21
22 # Step 4: Get best model performance from testing set
23 y_pred = best_clf.predict(X_ts_normalized)
24 test_acc = accuracy_score(y_ts, y_pred)
25 print(f'test_acc = {test_acc}')
```

```
gridsearch.best_params_ = {'C': 0.0001, 'penalty': 'l2'}
test_acc = 0.48
```

▼ RF

이번엔 random forest입니다.

```
1 # Step 2: Use GridSearchCV to find optimal hyperparameter values
2 clf = RandomForestClassifier()
3 parameters = {'n_estimators': [100, 150, 200],
4               'criterion': ['gini', 'entropy']}
5 gridsearch = GridSearchCV(clf, parameters, scoring='accuracy', cv=5)
6 gridsearch.fit(X_tr, y_tr)
7 print(f'gridsearch.best_params_ = {gridsearch.best_params_}')
8
9 # Step 3: Get model with best hyperparameters
10 best_clf = gridsearch.best_estimator_
11
12 # Step 4: Get best model performance from testing set
13 y_pred = best_clf.predict(X_ts)
14 test_acc = accuracy_score(y_ts, y_pred)
15 print(f'test_acc = {test_acc}')
```

```
gridsearch.best_params_ = {'criterion': 'gini', 'n_estimators': 100}
test_acc = 0.9966666666666667
```

Q: 두 모델의 예측 성능이 각각 어떻게 나타나나요? 주어진 데이터셋에 대하여, logistic regression과 random forest 중 어떤 모델이 더 적합할까요? 우리의 논의를 바탕으로, 이러한 성능 차이를 설명할 수 있나요?

References

- James Bourbeau. *Supervised machine learning in Python with scikit-learn*. URL: <https://jrbourbeau.github.io/madpy-ml-sklearn-2018/#/>
- Joaquin Vanschoren. *Machine Learning with Scikit-Learn*. URL: <http://joaquinvanschoren.github.io/ML-course-R/TutorialSKlearn.slides.html#/>
- User Guide - scikit-learn. URL: https://scikit-learn.org/stable/user_guide.html

