



1 사이킷런 소개와 특징

- 파이썬 머신러닝 라이브러리
- 특징
 - i) 쉬운 API
 - ii) 머신러닝을 위한 다양한 알고리즘과 개발을 위한 편리한 프레임워크/API 제공
 - iii) 활용도 ↑

2 머신러닝 - 분류

① 분류(classification)

- 지도학습 방법
 - ↳ 학습을 위한 다양한 feature와 분류 결정값인 label 데이터로 모델을 학습한 뒤, 별도의 test data 세트에서 미래의 레이블 예측

② Sklearn 내의 모듈

- i) sklearn.datasets : 사이킷런에서 자체적으로 제공하는 데이터 세트를 생성
- ii) sklearn.tree : 트리 기반 ML 알고리즘을 구현한 클래스의 모임
- iii) sklearn.model_selection : 데이터 분리(학습/검증/예측 데이터) or 최적의 하이퍼파라미터로 평가하기 위한 모듈

③ train - test - split () : dataset을 학습 데이터와 테스트 데이터로 분리

- ↳ parameters
 - feature data
 - label data
 - test_size : 전체 데이터 세트 중 test data의 비율
 - random_state : 랜덤 seed값 지정

↳ 반환값 > 순서대로 학습용 feature 데이터, 테스트용 feature 데이터, 학습용 label 데이터, 테스트용 label 데이터

④ model.fit(X_train, y_train) : 모델 학습시키기

분류 모델 (각각)

⑤ model.predict(X_test) : 테스트용 feature data에 대한 label값 예측

- ⑥ accuracy_score ()
 - 정확도 (예측 결과가 실제 label과 얼마나 일치하는가) 계산
 - accuracy_score (실제, 예측 결과)

3 사이킷런의 기반 framework

① Estimator 이해 & fit (), predict () 메서드

i) Estimator

- ↳ Classifier (분류 알고리즘을 구현한 클래스) + Regressor (회귀 알고리즘을 구현한 클래스)
- ↳ 학습: fit (), 예측: predict ()

② 주요 모듈

- ↳ 0.94 ~ 0.95 (표)

③ 내장된 데이터셋

i) 예제 데이터셋 API : datasets.load_데이터셋명 ()

ii) 표준 데이터 생성기

- a. datasets.make_classification () : 분류 데이터셋 생성
- b. datasets.make_blobs () : 클러스터링을 위한 데이터셋 생성

iii) key 값

- data : feature 데이터 세트
- target : 분류 ⇒ label, 회귀 ⇒ 숫자 결과값 데이터
- target_names : 개별 label의 이름
- feature_names : feature의 이름
- DESCR : 데이터 세트에 대한 설명 ⊕ 각 feature에 대한 요약



4 Model Selection 파트

- 학습 데이터와 테스트 데이터 세트 분리
- 교차 검증 분할·평가
- estimator의 하이퍼 파라미터를 튜닝

① train_test_split()

- 학습/테스트 데이터 분리
- ✱ 예측을 수행하는 데이터 세트는 학습을 수행한 학습용 데이터 세트가 아니라 전용 test 데이터셋이어야 함.
- parameters
 - i) feature datasets } 필수
 - ii) label datasets } 필수
 - *선택사항
 - iii) test_size : 전체 데이터 중 test data의 크기 지정 default: 0.25
 - iv) train_size : " train data "
 - v) shuffle : 데이터 분리 전 미리 섞을지 결정 default: True
 - vi) random_state : 시드값 지정 → 데이터 변경 방지
- 반환값
 - ↳ (X_train, X_test, y_train, y_test)
 - feature label

② 교차 검증

- 모델의 과적합(overfitting)을 방지
 - ↳ 특정한 학습 데이터에만 과도하게 최적화되는 것
- 데이터 편향을 막기 위해 별도의 여러 세트 구성된 학습 데이터 세트와 검증 데이터 세트에서 학습·평가 수행 → 파라미터 튜닝 등의 모델 최적화
 - 테스트 데이터로 최종 평가

i) k-fold 교차 검증

↳ k개의 데이터 폴드 세트를 만들어 k번만큼 각 fold 세트에 대해 학습-검증을 반복적으로 수행하는 방법

↳ 순서

1. 데이터 세트를 K등분
2. K번 반복하여 1개의 등분을 검증 데이터로, 나머지를 학습 데이터로 이용
 - ↳ 모든 등분이 한번씩은 검증 데이터로 사용되도록

↳ parameters & 메서드

i) kfold.split() } 데이터를 k개의 fold datasets로 분할
 feature data ↑ 리턴값 > 학습용/검증용 데이터로 분할할 수 있는 인덱스

ii) Stratified k-fold

↳ 불균형한 분포를 가진 레이블(정적 클래스) 데이터 집합을 위한 K-fold 방식
 ↳ 원본 데이터의 레이블 분포를 먼저 고려한 뒤, 이 분포와 동일하게 데이터 분배
 ↳ kfold.split(feature data, label data) 분포 고려

iii) cross_val_score()

↳ 교차 검증을 좀 더 편리하게 수행할 수 있는 API

↳ parameters

- 1) estimator : 사이킷런의 분류 알고리즘 클래스(분류/회귀) 적용 불가!
- 2) X : feature data
- 3) y : label data default > None
- 4) scoring : 예측 성능 평가 지표 " ex. accuracy : 정확도
- 5) cv : 교차 검증 fold 수 "

↳ 반환값 > scoring 파라미터로 지정된 성능 지표 측정값을 배열 형태로

↳ 내부적으로 Stratified K-fold를 이용

↳ 단 하나의 평가 지표만 지정할 수 0 cf > cross_validate()

↳ 여러 개의 평가 지표를 반환

↳ 학습 데이터에 대한 성능 평가·수행시간



② 피쳐 스케일링 & 정규화

- 피쳐 스케일링: 서로 다른 변수의 값 범위를 일정한 수준으로 맞추는 작업

i) 표준화

데이터의 feature 각각이 평균이 0이고 분산이 1인
가우시안 정규분포를 가진 값으로 변환

$$x_{i_new} = \frac{x_i - E(X)}{Sd(X)}$$

ii) 정규화

서로 다른 feature의 크기를 통일하기 위해 크기를 변환 \Rightarrow 단위 통일
사이킷런에서의 의미 > 개별 벡터의 크기를 맞추기 위해 변환

$$x_{i_new} = \frac{x_i}{\sqrt{x_i^2 + y_i^2 + z_i^2 + \dots}}$$

③ StandardScaler

- 표준화를 쉽게 자원하기 위한 클래스

- SVM, 선형회귀, 로지스틱 회귀에서 표준화가 중요

- scaler.fit() \rightarrow scaler.transform() \Rightarrow 표준화 결과가 np.ndarray로!

④ MinMaxScaler

- 데이터 값을 0과 1 사이 범위의 값으로 변환

- scaler.fit() \rightarrow scaler.transform() \Rightarrow 결과: np.ndarray

⑤ 스케일링 시 유의점

\hookrightarrow 학습 데이터로 fit()이 적용된 스케일링 기준 정보를 그대로 테스트 데이터에 적용해야 함.

\Rightarrow 데이터 스케일링 기준을 동일하게 적용!

\hookrightarrow 테스트 데이터에서 다시 fit()을 적용하면 X \Rightarrow 이미 fit()이 적용된 scaler 객체를 이용해 transform()을 적용해야 함.

\hookrightarrow 가능하다면 전체 데이터의 스케일링을 적용한 뒤 학습과 테스트 데이터로 분리하기