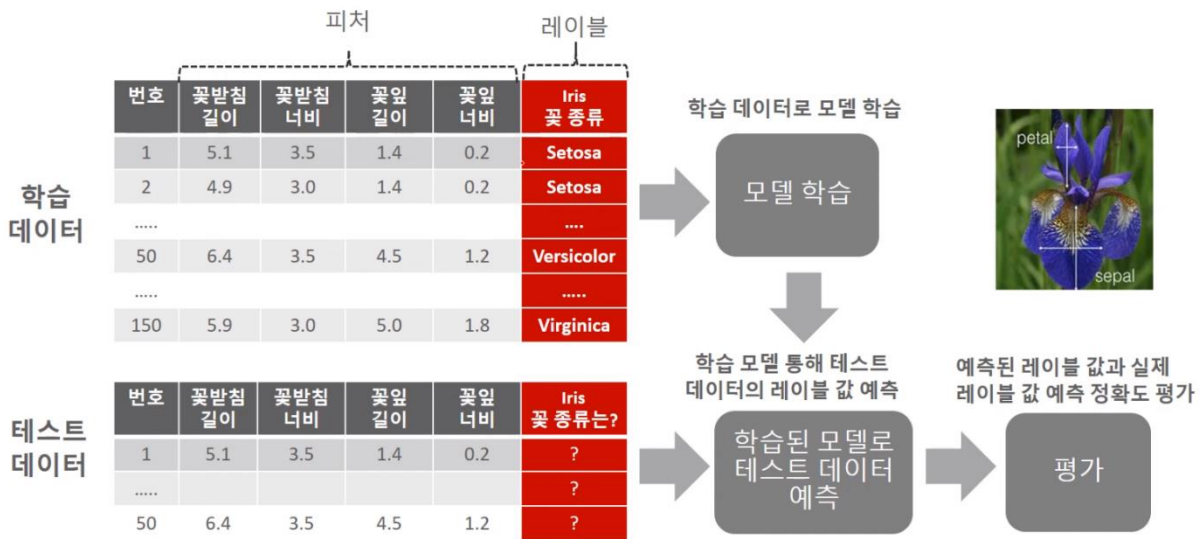
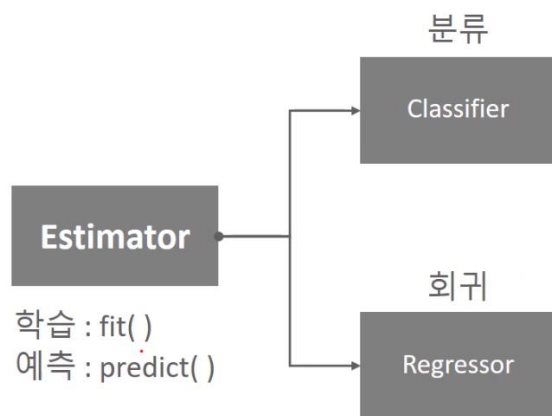


사이킷런의 기반 프레임 워크 익히기 - 주요 API/모듈 및 내장 예제 데이터 세트 소개

[붓꽃 데이터 분류 예측 프로세스]



[사이킷런 기반 프레임워크 - Estimator와 fit(), predict()]



분류 구현 클래스

DecisionTreeClassifier
RandomForestClassifier
GradientBoostingClassifier
GaussianNB
SVC

회귀 구현 클래스

LinearRegression
Ridge
Lasso
RandomForestRegressor
GradientBoostingRegressor

모든 구현 클래스를 통틀어서 사이킷런에서는 Estimator라고 부름

[사이킷런의 주요 모듈]

| 분류 | 모듈명 | 설명 |
|----------------------|----------------------------|---|
| 예제 데이터 | sklearn.datasets | 사이킷런에 내장되어 예제로 제공하는 데이터 세트 |
| 데이터 분리, 검증 & 파라미터 튜닝 | sklearn.model_selection | 교차 검증을 위한 학습용/테스트용 분리, 그리드 서치(Grid Search)로 최적 파라미터 추출 등의 API 제공 |
| 피처 처리 | sklearn.preprocessing | 데이터 전처리에 필요한 다양한 가공 기능 제공(문자열을 숫자형 코드 값으로 인코딩, 정규화, 스케일링 등) |
| | sklearn.feature_selection | 알고리즘에 큰 영향을 미치는 피처를 우선순위에 따라 선택 작업 수행하는 다양한 기능 제공 |
| | sklearn.feature_extraction | 텍스트 데이터나 이미지 데이터의 벡터화된 피처를 추출하는 데 사용됨. 예를 들어 텍스트 데이터에서 Count Vectorizer 나 Tfidf Vectorizer 등을 생성하는 기능 제공. 텍스트 데이터의 피처 추출은 sklearn.feature_extraction.text 모듈에, 이미지 데이터의 피처 추출은 sklearn.feature_extraction.image 모듈에 지원 API가 있음. |
| 피처 처리 & 차원 축소 | sklearn.decomposition | 차원 축소와 관련한 알고리즘을 지원하는 모듈임. PCA, NMF, Truncated SVD 등을 통해 차원 축소 기능을 수행할 수 있음 |

| 분류 | 모듈명 | 설명 |
|---------|----------------------|--|
| 평가 | sklearn.metrics | 분류, 회귀, 클러스터링, 페어와이즈(Pairwise)에 대한 다양한 성능 측정 방법 제공 Accuracy, Precision, Recall, ROC-AUC, RMSE 등 제공 |
| ML 알고리즘 | sklearn.ensemble | 앙상블 알고리즘 제공 랜덤 포레스트, 에이다 부스트, 그래디언트 부스팅 등을 제공 |
| | sklearn.linear_model | 주로 선형 회귀, 릿지(Ridge), 라쏘(Lasso) 및 로지스틱 회귀 등 회귀 관련 알고리즘을 지원. 또한 SGD(Stochastic Gradient Descent) 관련 알고리즘도 제공 |
| | sklearn.naive_bayes | 나이브 베이즈 알고리즘 제공. 가우시안 NB, 다항 분포 NB 등. |
| | sklearn.neighbors | 최근접 이웃 알고리즘 제공. K-NN 등 |
| | sklearn.svm | 서포트 벡터 머신 알고리즘 제공 |
| | sklearn.tree | 의사 결정 트리 알고리즘 제공 |
| | sklearn.cluster | 비지도 클러스터링 알고리즘 제공 (K-평균, 계층형, DBSCAN 등) |
| 유틸리티 | sklearn.pipeline | 피처 처리 등의 변환과 ML 알고리즘 학습, 예측 등을 함께 묶어서 실행할 수 있는 유틸리티 제공 |

[사이킷런 내장 예제 데이터 셋 - 분류 및 회귀용]

| API 명 | 설명 |
|--|---|
| <code>datasets.load_boston()</code> | 회귀 용도이며, 미국 보스턴의 집 피쳐들과 가격에 대한 데이터 세트 |
| <code>datasets.load_breast_cancer()</code> | 분류 용도이며, 위스콘신 유방암 피쳐들과 악성/음성 레이블 데이터 세트 |
| <code>datasets.load_diabetes()</code> | 회귀 용도이며, 당뇨 데이터 세트 |
| <code>datasets.load_digits()</code> | 분류 용도이며, 0에서 9까지 숫자의 이미지 픽셀 데이터 세트 |
| <code>datasets.load_iris()</code> | 분류 용도이며, 붓꽃에 대한 피쳐를 가진 데이터 세트 |

[내장 예제 데이터 셋 구성]

`feature_names`, `data`, `target_names`, `target`

| feature_names | | | | | target_names | | |
|---------------|----------------------|---------------------|----------------------|---------------------|-------------------------------|-------------|--|
| | sepal length (cm) | sepal width (cm) | petal length (cm) | petal width (cm) | setosa, versicolor, virginica | (0 , 1 , 2) | |
| data | 5.1 | 3.5 | 1.4 | 0.2 | 0 | target | |
| | 4.9 | 3.0 | 1.4 | 0.2 | 1 | | |
| | | | | | | | |
| | 4.6 | 3.1 | 1.5 | 0.2 | 2 | | |
| | 5.0 | 3.6 | 1.4 | 0.2 | 0 | | |

- 키는 보통 `data`, `target`, `target_name`, `feature_names`, `DESCR`로 구성
 - `data`: 피쳐의 데이터 세트
 - `target`: 분류-레이블 값, 회귀-숫자 결과값 데이터 세트
 - `target_names`: 개별 레이블의 이름
 - `feature_names`: 피쳐 이름
 - `DESCR`: 데이터 세트에 대한 설명과 각 피쳐의 설명