직장인을 위한 파이썬 데이터 분석

Scikit-learn Cheat Sheet

Scikit-learn



Scikit-learn 은 파이썬 대표 머신러닝 라이브러리입니다. 머신 러닝에서 사용될 수 있는 전처리, 학습, 예측, 평가 단계를 모두 제 공합니다.

라이브러리로딩

from sklearn.datasets import load_iris
from sklearn.preprocessing import StandardScaler
from neighbors import KNeighborsClassifier
from sklearn.cross_validation import train_test_split
from sklearn.metrics import accuracy_score

대이터로딩
전처리함수
학습모델
학습/테스트데이터분리
from sklearn.metrics import accuracy_score
평가지표

예제 사용법

iris = load iris() 데이터 로딩 X, Y데이터 지정 X, y = iris.data[:, :2], iris.target 학습 / 테스트 분할 X train, X test, y train, y test = train test split(X, y) scaler = StandardScaler().fit(X train) 표준화 X train = scaler.transform(X train) X test = scaler.transform(X test) knn = KNeighborsClassifier(n_neighbors=5) 학습 모델 생성 knn.fit(X train, y train) 모델 학습 y pred = knn.predict(X test) 모델 예측 평가 accuracy_score(y_test, y_pred)

학습 / 테스트 데이터 분할

import numpy as np
X = np.random.random((10,5))
y = np.array(['M','M','F','F','M','F','M','F','F','F'])
X[X < 0.7] = 0

from sklearn.cross_validation import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)

학습 모델 생성

지도 학습 (Supervised) 모델

from sklearn.linear_model import LinearRegression lr = LinearRegression(normalize=True)

from sklearn.svm import SVC 분류모델(SVM)
svc = SVC(kernel='linear')

from sklearn import neighbors 분류모델(KNN)
knn = neighbors.KNeighborsClassifier(n_neighbors=5)

비지도 학습 (Unsupervised) 모델

모델 학습 (X:문제, Y:정답)

모델 학습

지도 학습 (Supervised)

lr.fit(X, y)
knn.fit(X_train, y_train)
svc.fit(X_train, y_train)

비지도 학습 (Unsupervised)

k_means.fit(X_train) 모델 학습(X:문제) pca model = pca.fit transform(X train)

모델 예측

지도 예측 (Unsupervised)

y_pred = svc.predict(np.random.random((2,5))) 정답라벨(Label)예측
y_pred = lr.predict(X_test)
y_pred = knn.predict_proba(X_test)) 정답확률(Probability)예측

비지도 예측 (Unsupervised)

y pred = k means.predict(X test) 군집(Cluster) 예측

데이터 전처리

표준화 (Standardization)

from sklearn.preprocessing import StandardScaler 라이브러리로딩 scaler = StandardScaler().fit(X_train) 표준화스케일러 정의 standardized_X = scaler.transform(X_train) 표준화스케일러 적용

정규화 (Normalization)

from sklearn.preprocessing import Normalizer 라이브러리로딩 scaler = Normalizer().fit(X_train) 정규화스케일러정의 normalized X = scaler.transform(X train) 정규화스케일러적용

범주형 변수의 수치화 (Encoding)

from sklearn.preprocessing import LabelEncoder enc = LabelEncoder() 라벨인코더 정의 y = enc.fit_transform(y) 라벨인코더 적용

결측값 처리

from sklearn.preprocessing import Imputer 라이브러리로딩 imp = Imputer(missing_values=0, strategy='mean', axis=0) 대체기 (inputer) 정의 imp.fit transform(X train) 대체기 (inputer) 적용

모델 평가

분류 모델 평가

정확도 (Accuracy)

from sklearn.metrics import accuracy_score accuracy score(y test, y pred) 정확도 측정 (실제값 예측값 맞춘 비율)

오차 행렬 (Confusion Matrix)

from sklearn.metrics import confusion_matrix 오차행렬생성 print(confusion_matrix(y_test, y_pred)))

회귀 모델 평가

MAE (Mean Absolute Error)

from sklearn.metrics import mean_absolute_error 평균 절대 오차 y_true = [3, -0.5, 2] (평균적으로나타나는 오차) mean absolute error(y true, y pred))

MSE (Mean Squared Error)

from sklearn.metrics import mean_squared_error 평균제곱오차 mean_squared_error(y_test, y_pred)) (오차의넓이)

R² (R-Squared)

from sklearn.metrics import r2_score 적합도 r2 score(y true, y pred))

교차 검증 (Cross-Validation)

from sklearn.cross_validation import cross_val_score print(cross_val_score(knn, X_train, y_train, cv=4)) 분류교차검증(4분할) print(cross_val_score(lr, X, y, cv=2)) 회귀교차검증(2분할)

모델 튜닝

그리드 탐색 (GridSearchCV)

무작위 탐색 (RandomizedSearchCV)

from sklearn.grid search import RandomizedSearchCV 무작위 조합으로 탐색 params = {"n neighbors": range(1,5), 그리드 생성 "weights": ["uniform", "distance"]} rsearch = RandomizedSearchCV(estimator=knn. 무작위 탐색 모델 정의 param_distributions=params, 교차 검증 데이터 분할 수 n iter=8, 탐색 횟수 random state=5) rsearch.fit(X train, v train) 무작위 탐색 최적 파라미터 출력 print(rsearch.best params)