

새로운 특성(feature)를 선택하는 방법

학습 내용

- 01 일변량 통계(univariate statistics)
- 02 모델 기반 선택(model-based selection)
- 03 반복적 선택(iterative selection)

1-1-1 일변량 통계

- 개개의 특성과 타겟(목표변수) 사이에 중요한 통계적 관계가 있는지 계산
- 분류에서는 분산분석(ANOVA)라고 한다.
- 각 특성(feature)이 독립적으로 평가.
- 계산이 매우 빠르고 평가를 위한 모델을 만들 필요가 없음.
- SelectPercentile에서 특성을 선택하는 기준은 F-값. 값이 클수록 클래스 평균의 분산이 비교적 크다.

분류 - f_classif, 회귀 - f_regression

```
In [2]: import warnings
warnings.filterwarnings(action='ignore')
# warnings.filterwarnings(action='default')
```

```
In [3]: from sklearn.datasets import load_breast_cancer
from sklearn.feature_selection import SelectPercentile, f_classif
from sklearn.model_selection import train_test_split
import numpy as np
```

```
In [4]: cancer = load_breast_cancer()
print(cancer.data.shape)
```

(569, 30)

```
In [5]: # 고정된 난수를 발생
rng = np.random.RandomState(42)
noise = rng.normal(size=(len(cancer.data), 40))
noise.shape
```

Out[5]: (569, 40)

```
In [6]: # 데이터 노이즈 특성 추가
# 30개는 원본 특성, 다음 40개는 노이즈
X_w_noise = np.hstack([cancer.data, noise])
X_w_noise.shape
```

Out[6]: (569, 70)

```
In [7]: X = X_w_noise # 입력
y = cancer.target # 출력

X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                    random_state=0,
                                                    test_size=0.5)

# 50%를 뽑는 것을 학습
```

```
Out[7]: SelectPercentile(percentile=50)
```

```
X_train.shape: (284, 70)
X_train_selected.shape (284, 35)
```

- ```
In [9]: import matplotlib.pyplot as plt
```

A horizontal genomic map of chromosome 10, spanning from position 0 to 60 Mb. The map shows various bands of varying widths representing different genes or regions. A specific band labeled "HLA-DQA1" is highlighted in black, located between approximately 38 Mb and 42 Mb.

전체 특성 사용 : 0.940  
선택된 일부 특성 사용 : 0.923

## 2/4

- SelectFromModel은 지도학습 모델로 계산된 중요도가 임계치보다 큰 모든 특성을 선택
- 절반 가량의 특성이 선택될 수 있도록 중간값을 임계치로 사용.
- 트리 100개로 만든 랜덤 포레스트 분류기를 사용.

[illegible]

3/4

```

n_features_to_select=40)

select.fit(X_train, y_train)

선택된 특성을 표시합니다.
mask = select.get_support()
plt.matshow(mask.reshape(1,-1), cmap='gray_r')
plt.xlabel("특성 번호")

```

Out[16]: Text(0.5, 0, '특성 번호')



- 일변량 분석이나 모델 기반 특성보다 특성 선택이 나아짐.
- 랜덤 포레스트 모델은 특성이 누락될때마다 다시 학습하므로 40번 실행.
- 이 코드를 실행하면 모델 기반 선택보다 훨씬 오래 걸림.

```

In [17]: X_train_rfe = select.transform(X_train)
X_test_rfe = select.transform(X_test)

score = LogisticRegression().fit(X_train_rfe, y_train).score(X_test_rfe, y_test)
print("테스트 점수 : {:.3f}".format(score))

```

테스트 점수 : 0.923

```

In [18]: ### RFE에서 사용된 모델로 예측
print("테스트 점수 : {:.3f}".format(select.score(X_test, y_test)))

```

테스트 점수 : 0.933

In [ ]: