

Chapter 03. 피쳐 엔지니어링과 모델 규제. 데이터의 무게를 예측하라!

03-1. K-최근접 이웃 평가.

문제가 되는 문제가 아니고 모델을 예측해야 한다. → 무게는 어떤 숫자 가능.

→ "회귀"

→ 무거운 것일수록 샘플을 ~~중요~~. → 데이터 샘플을 ~~가중~~ get!

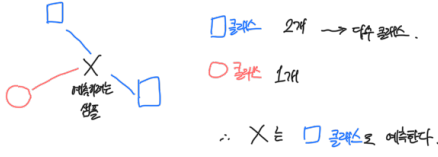
지도학습 알고리즘

분류 : 샘플을 몇 개의 클래스 중 하나로 분류하는 문제. → 정답이 있음.

회귀 : 임의의 입력 수를 예측하는 문제. → 두 배의 데이터 샘플을 가중치를 부여하는 방법.
→ 데이터의 무게 예측.

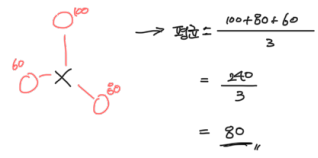
K-최근접 이웃 분류 알고리즘

- 1. 예측하려는 샘플에 가장 가까운 샘플 K개를 선택.
- 2. 샘플들의 클래스를 확인하여 다수 클래스를 새로운 샘플의 클래스로 예측한다.



K-최근접 회귀 알고리즘

- 1. 예측하려는 샘플에 가장 가까운 샘플 K개 선택.
- 2. 샘플들의 평균을 구한다.



<데이터 준비>

- 데이터셋 : 데이터의 집합. 길이 : 특성, 무게 : 타겟.

- 샘플링 : 데이터의 일부를 추출하는 과정. → 데이터셋 import 후 subset() 함수 사용해서 샘플링 가능.

사이킷런의 train_test_split() 함수 : 훈련 세트, 테스트 세트 나열.
2개의 배열 샘플링 가능.
→ 1차원 배열 → 2차원 배열. [1, 2, 3] → [[1], [2], [3]]
크기: (3, 1) → 2차원 배열 만들기 위해 크기로 바꿔 줄 필요가 있다. → 배열을 나타내는 불변형 배열, numpy 배열.

특성을 하나만 사용 → 샘플을 2차원 배열 만들어야 함. → 크기가 바뀐 크기를 바꿀 수 있는 reshape() 메소드 사용.

* reshape() 메소드

: 배열의 크기 바꿈.

→ 바꿀 때는 배열의 크기를 지정할 수 있다.

```
ex) test_array = np.array([1, 2, 3, 4]) // (4, ) 배열.  
test_array = test_array.reshape(2, 2) // (2, 2) 크기로 바꿈.  
print(test_array) // (2, 2) 출력.
```

* 넘피이는 배열의 크기를 자유롭게 지정하는 기능도 제공한다.

크기 -1 : 나머지 요소 개수로 모두 채우는 의미.

```
ex) train_input = train_input.reshape(-1, 1)  
test_input = test_input.reshape(-1, 1)  
print(train_input, test_input)
```

reshape(-1, 1) → 배열의 정사각형 구조 필요 X.

from sklearn.neighbors import KNeighborsRegressor k-최근접 이웃 회귀 알고리즘을 구현할 클래스 KNeighborsRegressor.

knr = KNeighborsRegressor() **가져오기**

knr.fit(train_input, test_target) **k-최근접 이웃 회귀 모델 훈련.**

print(knr.score(test_input, test_target))

⇒ 0.9928094061010639

[분류 : 정확도 (정답에 대한 개수 비율)
회귀 : 결정계수 (R²)

* $R^2 = 1 - \frac{(\text{예}-\text{실})^2 \text{의 합}}{(\text{예}-\text{평균})^2 \text{의 합}}$

↳ 타깃의 평균 정도를 예측하는 쉬운 (빠, 분포가 비슷함) → R²은 0에 가까워짐.

↳ 예측이 타겟에 아주 가까워지면 (숫자가 0에 가까워짐) → R²은 1에 가까워짐.

* **mean_absolute_error** : 타깃과 예측의 절대값 오차를 평균하여 반환함.

훈련세트와 테스트 세트의 결과를 비교했을 때

과대적합 : 훈련세트가 너무 높으면 "과대적합"

↳ 훈련세트에서 점수가 굉장히 좋았는데 테스트 세트에서는 점수가 굉장히 나쁘다. ⇒ 훈련세트에 과대적합이다.

⇒ 훈련 세트에만 잘 맞는 모델. ↳ 원치 X.

과소적합 : 2 변수거나 수가 모두 낮으면 "과소적합"

↳ 훈련세트보다 테스트 세트의 점수가 높거나 두 점수가 모두 너무 낮다. ⇒ 훈련세트에 과소적합이다.

⇒ 모델이 너무 단순하여 훈련 세트에 적절히 훈련되지 않은 경우.

↓
과소적합 해결방법 : 모델을 더 복잡하게 생성.

↓
k-최근접 이웃 알고리즘 모델에서 더 복잡하게 만드는 방법 : 이웃의 개수 k를 늘림.

* $\left[\begin{array}{l} \text{이웃의 개수} \downarrow : \text{훈련 세트에 있는 유사적인 패턴에 민감해짐.} \\ \text{이웃의 개수} \uparrow : \text{데이터 전반에 있는 일반적인 패턴을 따름.} \end{array} \right.$

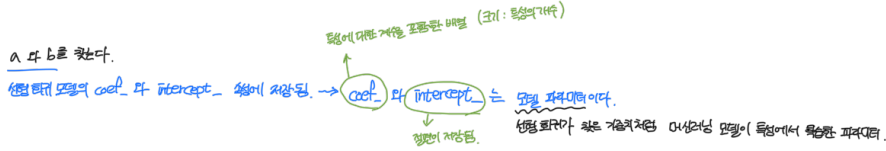
선형회귀 - 널리 사용되는 대표적인 회귀 알고리즘.

: 특징이 특정한 경우 어떤 값을 예측하는 알고리즘.
↳ 특정 값을 가질 수 있는 조건.

* 파이썬은 sklearn.linear_model 패키지 아래 LinearRegression 클래스로 선형 회귀 알고리즘 구현.

↳ 가장 간단한 회귀 모형의 사용.

↳ LinearRegression 클래스가 이 데이터에 가장 잘 맞는 a 와 b 를 찾는다.



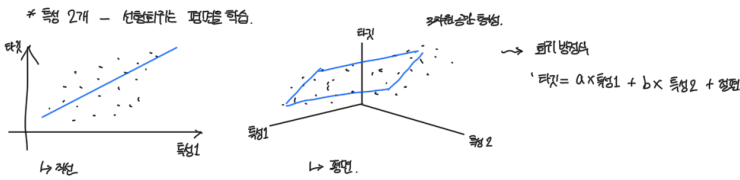
다항회귀 : 다항을 회귀한 선형 회귀

선형도 - 곡선 → 가장 적절한 곡선을 찾기. (여러개의 그래프와 가장 적절한 곡선이 특정에서 특정한 파라미터를 추출해야 한다.)

2차방정식 → 선형회귀 진행.

03-3. 특성 공학과 규제

다중회귀 : 여러 개의 특성을 사용한 선형회귀



* 특성 3개 - ✕

→ 3차원 공간 이상을 그려거나 상상 ✕

특성공약 : 가변의 품질을 사용해 새로운 특성을 뽑아내는 작업.

* 판다스 : 데이터 분석 라이브러리.

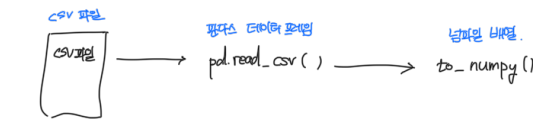
* 데이터 프레임 : 판다스의 핵심 데이터 구조.

→ 데이터 배열과 비슷하게 다차원 배열을 더할 수 있지만 행선 더 많은 기능 제공.

판다스에서 파일 읽는 법

: `read_csv()` 함수에 주소를 넣음 → `pd.read_csv(주소)` 데이터 프레임 만들기.

`to_numpy()` 메소드 사용해서 넘파이 배열로 변경. ex) `perch_full = df.to_numpy()`



변환기 : 특성을 만들거나 강력하게 유해 제공하는 다양한 클래스

→ 사이킷런에는 변환기라 부르는 특성을 만들거나 강력하게 제공하는 클래스가 있다.

→ 변환기 클래스는 `fit()`, `transform()` 메소드를 제공.

훈련 변환 → 학습 데이터에 `fit_transform` 메소드 존재.

PolynomialFeatures 클래스

→ `sklearn.preprocessing` 패키지에 포함.

규제 : 머신러닝 모델이 훈련 세트를 너무 과도하게 학습하지 못하도록 제한하는 것.

→ 모델이 훈련세트에 과대적합되지 않도록 만들.

→ 훈련된 모델의 경우 특성에 공약하는 계수(아 가중치)의 크기를 작게 만들.

* 선형회귀 모델에 규제를 추가한 모델을 '릿지'와 '라소'라고 부른다.

- 릿지 : 계수를 제공한 값을 기준으로 규제 적용, 선형모델의 계수를 작게 만들어 과대적합을 방지시킨다. (비록 효과가 좋아 보일 수 있음).
- 라소 : 계수의 절대값을 기준으로 규제 적용, 릿지다 달리 계수 값을 아예 0으로 만들 수도 있다.

- 릿지 회귀 : 모든 계수를 만들고 `fit()` 메소드에서 훈련한 뒤 `score()` 메소드로 평가한다.
- 라소 회귀 : `Bridge` 클래스 → `Lasso` 클래스 변경.