

기계학

김문 / 응답

202184026

서종현

201904025

김성욱

데이터를 나누는
다른 방법이 있는가?

①

서용현: data를 임의적으로 하기위해
train-test, test-set을 7:3 으로 나누면
어느 상황이든 임의적으로 할 수 있는 것 같다.
7:3

김성욱: 보통 data를 train-set, test-set
을 8:2로 나누는데 좀더 정확한 결과를
하기위해 Train-set의 8에서
Train-test 할부분을 나눈다.

$\therefore (\text{Train-data}, \text{Train-test}) : \text{Test-data}$
 $= 8:2$

② 무작위 샘플링 VS 계층적 샘플링 ?

서용현: 계층적 샘플링이 좋은 것 같다.
표본량 불균형이 적어 더 좋은 예측을
할 수 있는 것 같다.

김성욱: 무작위 샘플링도 편향을 줄 수 있는게

정해진 틀에서 데이터를 나누면 무작위적인
데이터가 들어왔을때 오해시켜 예측을
못 할 것 같다.

지도?
경험?

김성욱: 지도, 비지도, 강화학습은
고르는 기준이 뭐가?
무엇보고 결정하리?

서종현 = 지도는 경험, 비지도는
경험 무, 강화학습은 보정이
필요할 때

김성욱: Label 이 무조건 이루어져있나?

서종현 = > 예

서종현: 데이터 셋에서 null로
채워지면 null값을 왜
처리해야 하는가?

김성욱 = null값으로 채워진 데이터는
학습에 아무런 영향도 주지
않는 데이터이기 때문에 학습에
도움주는 방법으로

바뀌어주는게 학습한 모델의
정확도를 올려준다.

더 오래 기억하기 위한 부분

**** 전처리 ****

데이터 전처리 : 모델 학습을 효율적으로 진행하기 위해 주어진 데이터를 변환시키는 것

1. 수치형 데이터 전처리 과정 :

*null 값 처리 방법

- 1) 해당 구역 제거 - dropna
- 2) 전체 특성 제거 - drop
- 3) 특정 값으로 채우기 - 0, 평균, 중간값 등

2. 텍스트와 범주형 :

* 원-핫 인코딩 : 더미 특성을 추가하여 활용

- 1) 해당 카테고리의 특성 값 : 1
 - 2) 나머지 카테고리의 특성 값 : 0
- OneHotEncoder 클래스 제공

3. 나만의 변환

4. 특성 스케일 : 수치 데이터 전처리 과정