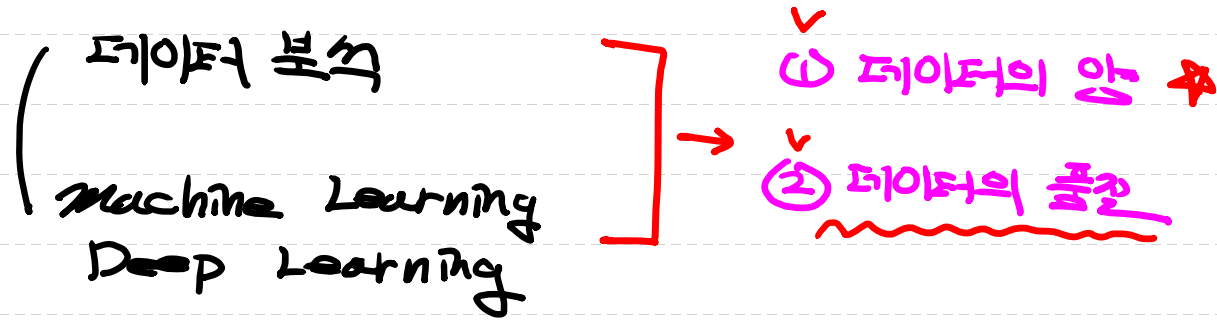


• 03/18 데이터 전처리



- ★ → Missing Value (결측치)
- ★ → Outlier (이상치)
- ★ → 인코더
- ★ → Data의 전처리
- ★ → 처리 → 1억 ~ 5억 ✓ → 10000000 ~ 50000000 ✓
● 1년 40년 ✓ → 1 ~ 40 ?

“Titanic” data Set

★ Seaborn module이 제공하는
Titanic data set을 이용

① Missing Value (결치값)

DataFrame 안에 **누락**된 값이 존재. (데이터 없거나 손실로 누락)
(파일 보거나 처리 문제로 누락)

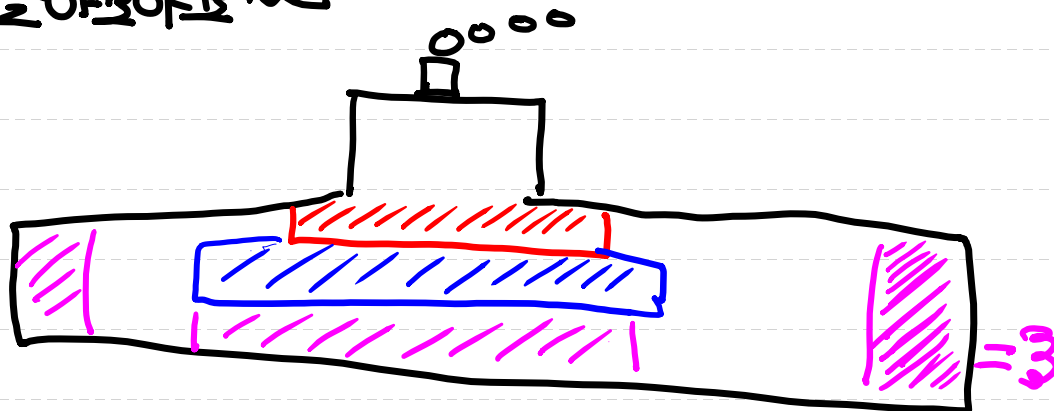
↓
NaN으로 표현

이런 Missing Value가 많으면 좋지 않아요!

처리해야 해요 ⇒ ① 삭제 → 만약 데이터가 충분히 많다면 그리고
Missing Value가 전체의 3~4% 이하일때

② 대체 → 평균, 중앙, 최대, 최소, 빈도
* 머신러닝 기법으로 이 값을 예측해서 채워요 *
↳ 충분한 자료를 안되고 조건이 있어요!

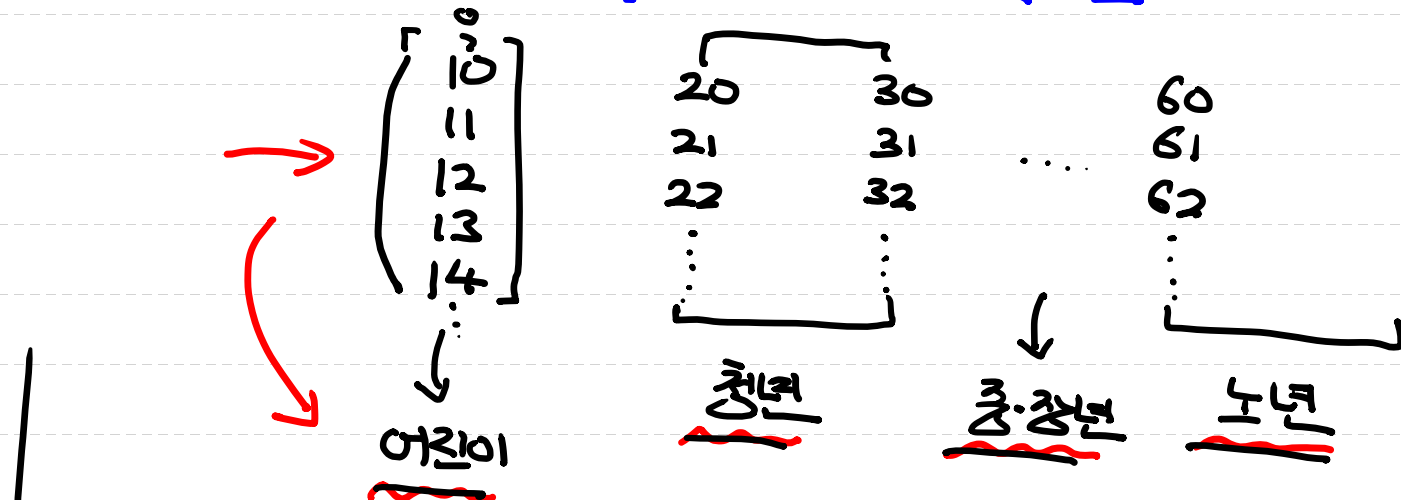
코드로 알아보세요~



② 이상치 처리 ~ "이상치를 찾아서 구현" ~ 나중에 다시 설명!

③ 중복 처리 ~ duplicated <
 drop_duplicates >

* ④ 자료형 브린 ~ * 만약 숫자가 문자열 (pandas Object)로 저장되어 있으면
숫자형으로 브린해주시기 맞아요!



→ Code로 살펴보아요~

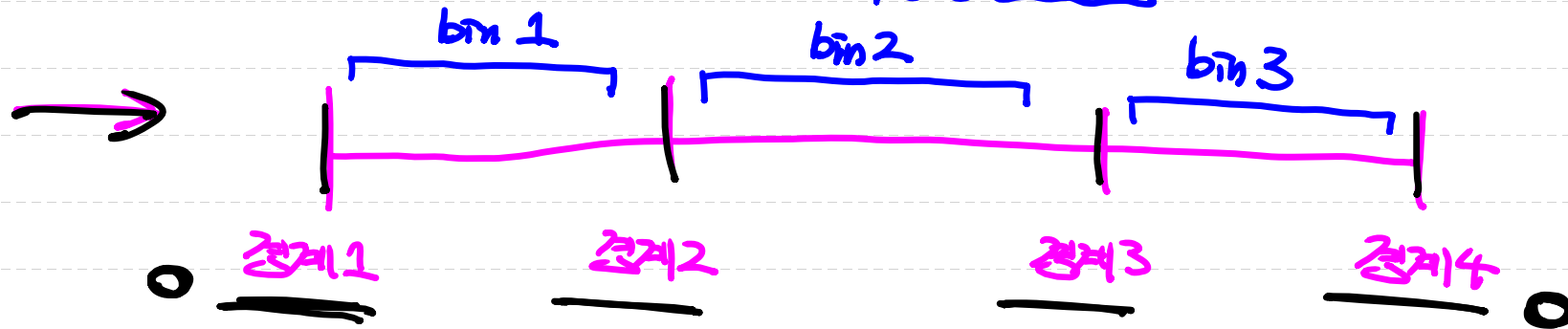
범주형 (Category) 데이터 처리

0세 ~ 15세 : 어린이
16세 ~ 25세 : 청소년
26세 ~ 40세 : 청년
...

→ 이런 형태로 연속적인 데이터를 범주형 데이터로
처리하는데 훨씬 효율적인 경우가 있어요!

* 구간분할

↳ 연속적인 데이터를 일정한 구간으로 나누기.



↓
코드로 구간분할을 통해 범주형 데이터를 만들어 보아요 ~

* housepower를 구간분할

- 저렴한
- 보통가격
- 고급

→ 범주형 데이터로 변환

Category를 나타내는 범주형 data는 Machine Learning 알고리즘에서

바로 사용하기 힘들어요!! → 컴퓨터가 인식할 수 있는 형태로 제공되어야
해요!

↓
이 문제를 해결하기 위해 dummy variable (더미 변수)

↓
0과 1로 표현하고 해당 특성이 있는지 여부를 표현

★ "One-hot-encoding"

⑤ Normalization (정규화)

DataFrame의 각 열에 가지고 있는 숫자 데이터의 상대적인 차이 때문에

머신러닝결과가 달라질수 있어요!!

→ 숫자 데이터의 상대적인 크기 차이를 제거/해결해야해요!!

*** <~

① 표준화 (Standardization) → "표준화된 z-score" (구분해 줘요)

크약점이
올까요??

② Min-Max Scaling → 알아보아요 ~

★

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

$$\frac{3 - 1}{5 - 1} = \frac{2}{4}$$

	A	B	C
0	1000	2	20
1	3000	1	50
2	5000	3	70
3	2000	5	10000

	A	B
0	0.25	
0.5	0	
1	0.5	
0.25	1	

이상이기
조금하면
이상해줘요