

# Definitive EDA

2023.07.01 Sat  
DATE NO.

distplot → 연속형 값. #152행.

distribution(분포) 확인은 중요해.

>> 13 numeric feature와 비교가능하다.  
alphabet과 string은 비교X

피어슨 상관계수

: 각 feature의 분산들이 어떤식으로 되어있어서 어떤 상관관계가 있나?  
어떤 형태의 모양을 보고 그것을 수치화해서 그들의 correlation을 본다.

다중공선성  $y=x$  .. 똑같은 값을 가질 때. 그러다 redundant해.

필요없어. 비슷한 값을 가진것이 많으면 성능저하

예) 항상불같은 부스럼 모델은 여러 트리를 합쳐서 만드는 모델

이러한 모델의 가치는 다양한 (diversity)를 가진 트리를 만들어야

성능이 좋아진다. EDA를, feature engineering and data cleaning.

모델 학습은 컴퓨터가 알아서 함. 사람이 해야하는건 설계, EDA, idea, 평가하는 문제! 지이러 이해가 커진다.

모델이 학습하기 좋은 형태로의 변환도 중요하다.

>>

나쁜것은 연속형 값이다. 어느문제가있다.

13개의 경우 남/여로 (Sex)로 나눌수있지만 Age로 나눌수  
없.. 나눌수없다. 문제가있다. 나뉘어야. 그래서 카테고리값으로..

마를 필요가있어. range-hand로 풀자.

80까지 max까지 5개 bins.  $80-16=5$ .

bins-size = 16.



→ 연속적인 value와 이산적인 value를 비교하면 어떻게 유무를 잘 판별 가능  
Survived 같은 binary value는 비교가 편하다.

Let  $f$  be continuous - Value  $\frac{2}{2}$  is achieved.

특정 sample 42 보러간다. range안에서

^ Out으로 나누면 같은 것임 보인다. (구간의 총)

\*  $\rho_{\text{air}}$ 은 같은 ~~조건~~ 수로...  $\rho_{\text{air}}$ 을 쓰라.  $\rho_{\text{air}}$ 을  $\rho_{\text{air}}$ 로 판독하라.

cut,  $g_{cut}$  둘 다 확인하자

» 각변 인코딩을 활용해서 String → 숫자로.

data feature from factor into?

Style

⇒ h2.1.2가 innermost와 style을 조합하여 만드는 HTML Element  
 55페이지 h2.3 innermost와 style을의 조합에이 접어다.

Document의 장으로 구성되어 있다.

document # 21841 same

document, dir1/101, dir2/102, ...

12/27/94

document. QuerySelector('h1')

let  $\alpha = \frac{1}{2}$

```
fontObject.setFontText = '1024px 40px';
```