# DTE Tanic EDA

data.pclass, data.survived 이 방식보다
data["Sex"], data["pclass"] 추천합니다.
가시성이 훨씬 좋아요.

>> Crosstab → sns.factorplot 으로 가시화
(ex group by랑 비슷.

이거 안에서 → catplot



error bar.
표준편차

Pclass feature는 중요하거나
유의미해!

target과의 연관성.

>>
Age같은 연속형 값은 히스토램을 많이 그린다.
왜냐하면 함수형태인 kdeplot을 이용하자. 가우시안 커브.

Null값을 해결할때 group화을 해서 평균을 구해 넣어주면 좀 낫다.

feature engineering 은 데이터을 변형해서 insight를 얻거나
우리가 가진 상상을 가지고 대처할 상현다음에 그것을 모델에 넣어주고
학습을 해서 더좋은 성능을 만드것.

>>
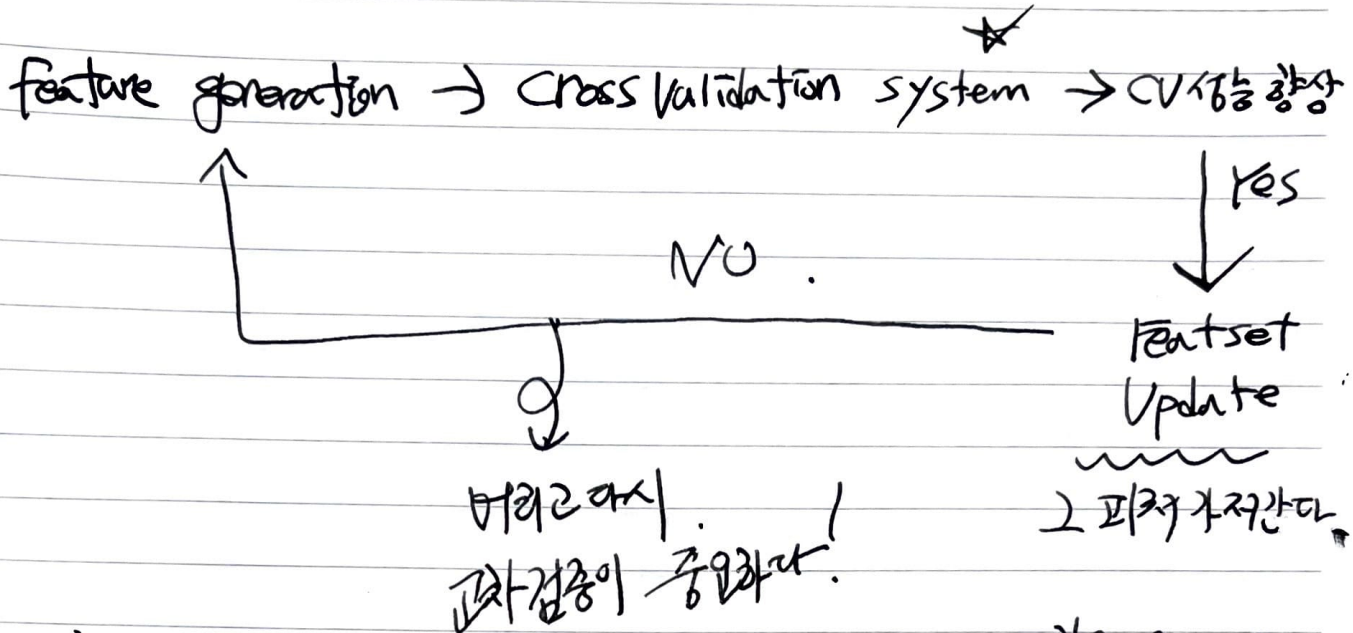feature을 변경해주었을 때 주의할 것은 'inplace = True'를 지정해서
원래값이 없는 방법이 있고, 다른 방법으로는 변수안에 넣는것이다.

>>
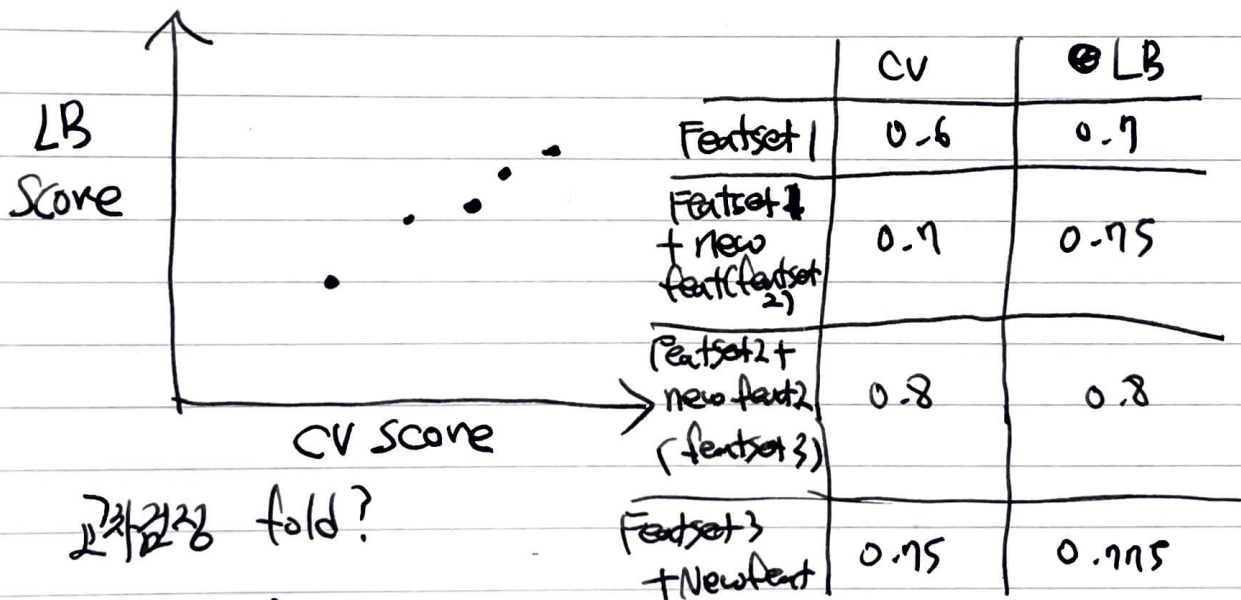data.loc [(data.Age.isnull( )), : ]
기 location 안에 mask (조건)을 넣으면 참 편리하다.

1. data .isnull().sum()
2. EDA
3. 교차검증 System.
   (Cross validation)

feature generation → Cross validation system → CV성능 향상

↑                                          ↓ Yes

                    NO.

                                          Featset
                                          Update

        ♀                                 그 피처 가져간다.

    머리박시.
    교차검증이 중요하다!

>> Feature generation과 Cross-validation은 함께 해야 된다.



| | CV | ⊛ LB |
|---|---|---|
| Featset 1 | 0.6 | 0.7 |
| Featset 1 + new feat(featset 2) | 0.7 | 0.75 |
| Featset2+ new feat2 (featset 3) | 0.8 | 0.8 |
| Featset3 + Newfeat | 0.75 | 0.775 |

LB Score

CV Score

교차검증 fold?

K-fold
LOO
LPO
:

→ CV와 LB가 같이 움직인다.
  좋은 교차검증 중에는 좋은 feat이 있다.

  CV와 LB가 비슷하게 찾올려가는
  교차검증 System 만들자.