

Data Preprocessing (데이터 전처리)

★ ① 다양한 양질의 데이터 → 필수처리
아래처리
정규화 처리 (Normalization)

★ ② Model을 만듭니다.
 { Python ~> 속도가 느려요 (x)
 Sklearn ~> 권장됨!! (o)
 tensorflow
 pytorch } • 1.x 버전으로 공부
 2.x 구현할때.

• Min-Max Scaling (이항치에 민감) → 0 ~ 1 사이의 값으로
 • Standardization (표준화) → 평균과 표준편차로
 Student's T 분포
 Scaling은 중요합니다.
 (실제로서
 이항치에 민감)

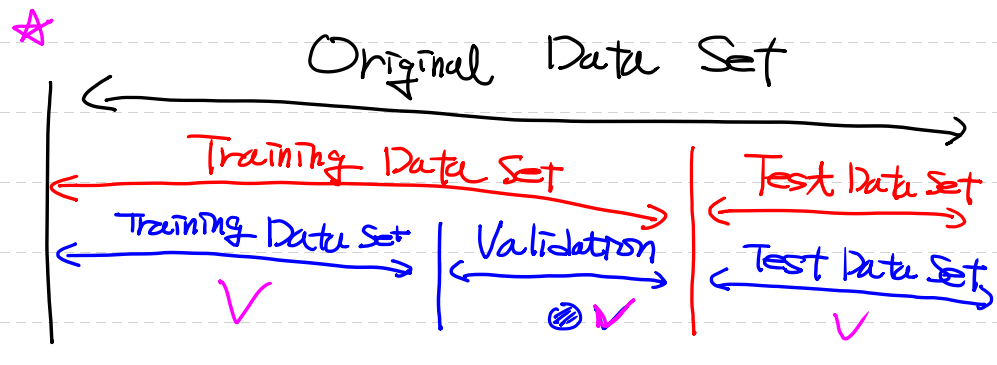
★ ③ Evaluation (평가) : Metric

- ① precision
- ② recall
- ③ *** accuracy
- ④ F1
- ⑤ Fall-out

• Training Data Set을 이용해서 Accuracy를 측정하면
 안되요!



• ~~선형~~가
Data Set



Training Data Set: 학습용도

Test Data Set: 최종 성능 평가.

Validation Data Set: 모델 개선 용도의
Test용
data set

• Data Set이 충분히 많지 않을 경우 ☹

→ Training Data Set 작을수록 성능이 낮아져요

(과소적합)

underfitting

epoch을 늘려서 학습

"k-fold Cross Validation" (CV)

"교차검증"

우리 Model이
해당 data만

를 표현하고

일반적인 data에

대응하는 것

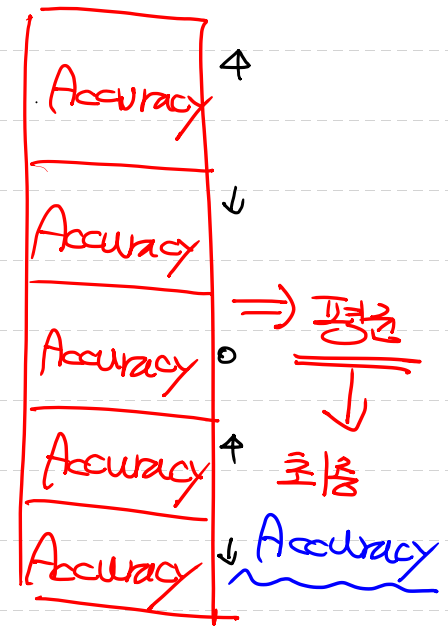
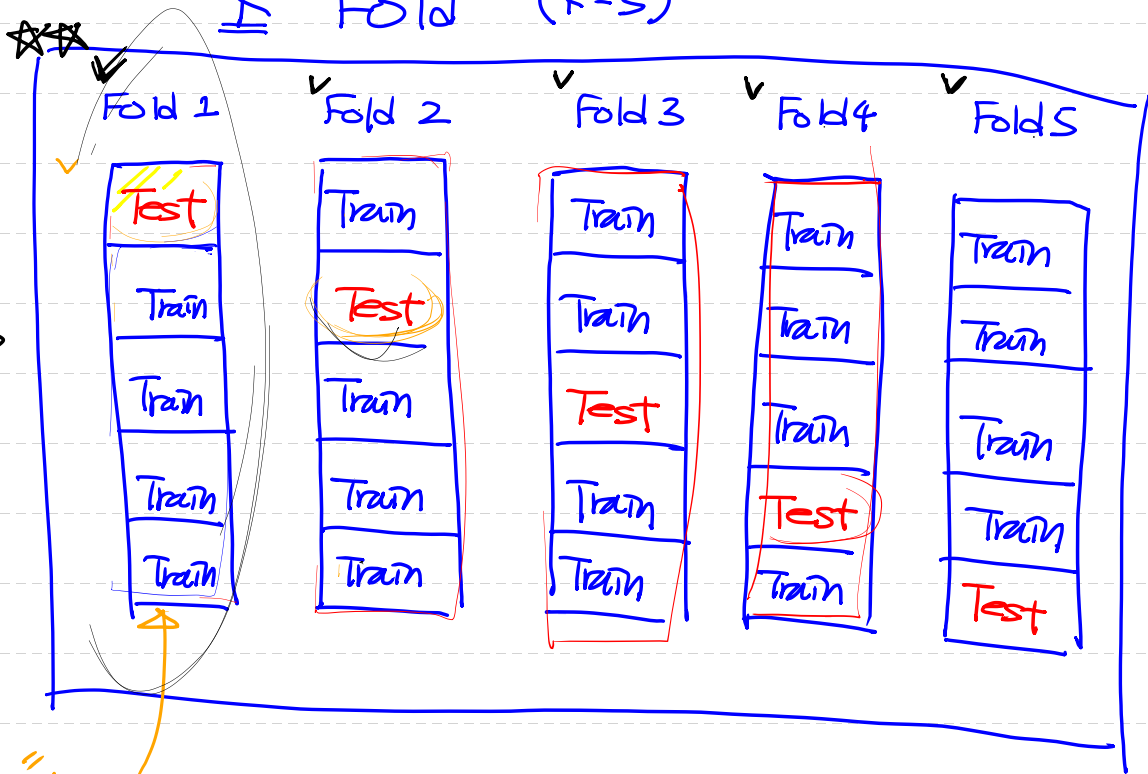
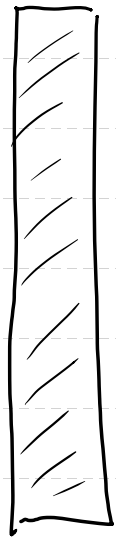
안되는 현상

(과대적합)

→ overfitting

"K-Fold" (k=5)

"Data"



구현 !! → "BMI 예제". 선처리
모델생성
accuracy