

## 머신러닝 과정

데이터 가공	1. 데이터 클리닝	데이터 분석
	2. 데이터 확인 (시각화)	
	3. 데이터 전처리	
모델 학습 예측	4. 머신러닝 알고리즘에 적합한 데이터 준비	scikit-learn tensorflow pytorch
	5. 머신러닝 알고리즘 적용	
	6. 머신러닝 알고리즘 평가	
평가 반복	7. 머신러닝 알고리즘 재현	
	8. 결과 나올 때까지 1~7 반복	

### 회귀와 분류

- 평가 지표가 다르다.
- 결과 차이 분명

회귀	RMSE	RMSLE
분류	F1 score	

### 회귀 평가 준비

MAE	Mean Absolute	$\frac{1}{N} \sum_{i=1}^N  y_i - \hat{y}_i $	실제, 예측 값 차이 절댓값 평균 $\Rightarrow$ 평균
MSE	Mean Squared	$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2$	차이 제곱 평균
RMSE	Root MS	$\sqrt{\frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2}$	$\sqrt{MSE}$ MSE는 너무 큼.

MSLE	Mean squared log	$\frac{1}{N} \sum_{i=1}^N (\log(y_i+1) - \log(\hat{y}_i+1))^2$	MSE에 log 적용 $\log(y+1)$
RMSLE	Root MLS	$\sqrt{\frac{1}{N} \sum_{i=1}^N (\log(y_i+1) - \log(\hat{y}_i+1))^2}$	$\sqrt{MSLE}$
R <sup>2</sup>	R Square	예측값 분산 / 실제값 분산	

### 분류 평가 지표

- 정확도 :  $\frac{\text{예측 결과 중 실 데이터}}{\text{전체 예측 데이터}}$  직관적, label 불균형한 경우 성능 더욱 저하

오차 행렬 (혼동 행렬) confusion matrix

		예측 클래스 $\hat{y}$	
		Negative 0	Positive 1
실제 클래스 y	Negative 0	TN F $\rightarrow$ F	FP F $\rightarrow$ T
	Positive 1	FN T $\rightarrow$ F	TP T $\rightarrow$ T

$\Rightarrow$  예측 클래스에 대한 평가

① 합치나 거짓이나?

② 예측을 어떻게 했나?

$$* \text{오차 행렬에서의 정확도} = \frac{TN + TP}{TN + FP + FN + TP}$$

- 정밀도와 재현율 : positive 데이터에 초점

- 정밀도  $\frac{TP}{FP + TP}$  양성 예측도

모델이 true로 분류한 데이터 중 실제 true 비율.

ex) 예측한 날씨가 실제로 얼마나 맞았나?

- 재현율  $\frac{TP}{FN + TP}$  민감도

실제 true인 데이터 중 모델이 true로 분류한 데이터 비율.

ex) 실제 날씨가 맑은 날에 대해 모델이 맑다고 예측한 비율.

\* 서로 연관 관계.

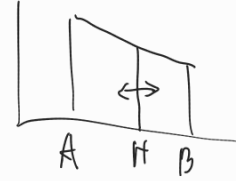
- F1 score : 정밀도, 재현율 결합한 지표.

$$F1 = \frac{2}{\frac{1}{\text{recall}} + \frac{1}{\text{precision}}} = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

(조화 평균)

산술 평균에 비해 bias가 ↓

⇒ label이 불균형한 경우에도 적용 가능.



· ROC 곡선과 AUC

이진 분류 모델 성능 측정

$$- FFR = \frac{FP}{FP + TN}$$

