

< 간단한 머신러닝 >

Date.

No.

Chapter 02. 모델 평가 및 선택

2.1 경험 과다 및 과적합

m 개 샘플, a 개 샘플이 오분류 ($m > a$)

$$E \text{ (error rate)} = \frac{a}{m}, \quad \text{Accuracy (정확도)} = 1 - \frac{a}{m}$$

$\left\{ \begin{array}{l} \text{Training error / Empirical error} \\ \text{★ Generalization error} \end{array} \right.$

"새로운 샘플 데이터에 대해 그 좋은 성능을 반복하는 학습기" → 원의 목표!

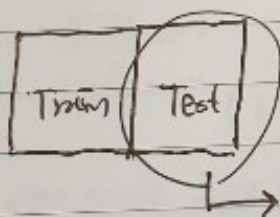
↳ 이를 위해 (Train) 데이터에서 (loss)의 잠재적인 "편향성"을 찾아야 함!

2.2.4 $\left\{ \begin{array}{l} \text{학습이 과다하기 같다면} \rightarrow \text{일반화 성능} \downarrow \text{ "Overfitting"} \\ \text{학습이 과소하기 같다면} \rightarrow \text{학습도 제대로 X "Underfitting"} \end{array} \right.$

Then, 현실 프로젝트에 직면하여, ① 어떤 학습 알고리즘을 사용?

② 어떠한 파라미터를 선택?

→ Model Selection



→ 학습기의 일반화 능력에 대해 평가!! ★ train과 겹치면 X

Test error $\overset{\text{Approx.}}{\approx}$ Generalization error.

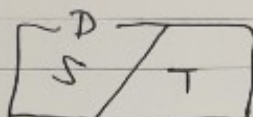
Date. 2.2 평가 방법

Assume that,

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$$

Data set D 를 어떻게 사용할 것인가? 훈련셋 S 와 테스트셋 T 로 나눌 것인가?

① Hold-out



안으로 잘라 놓게..!

* 데이터의 분포가 같게 나뉘어야 함! → 데이터 분포를 유지하는 결과를 얻을 수 있음.

예, 층화 추출법 (Stratified Sampling) 등 고려.

→ 층화 추출법으로 문제가 존재한다고 함.

(클래스 불균형) 클래스 사이 이질성을 갖춰야 할 문제들에서 고려하는 것이 어려움!

* 안으로 나뉘지 않으면 어떻게 추출해야 하냐 (train 9)

평가 결과가 달라짐.

+ 비율이 아슬아슬 존재! (일반적으로 $\frac{2}{3} \sim \frac{1}{3}$)

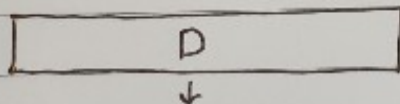
② Cross Validation

$$D = \bigcup_{i=1}^k D_i, \quad D_i \cap D_j = \emptyset \quad (i \neq j)$$

→ 최대한 데이터 분포가 나뉘게!

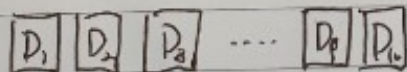
층화 추출!!

k-fold cross validation. (일반적으로 $k=10$)

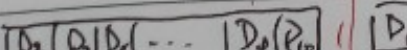
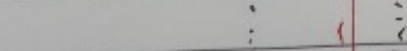
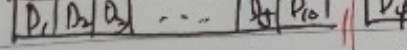
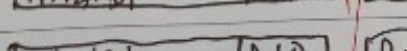
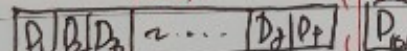


+ LOOCV (Leave-One-Out CV)

$k=m$ where $|D|=m$.

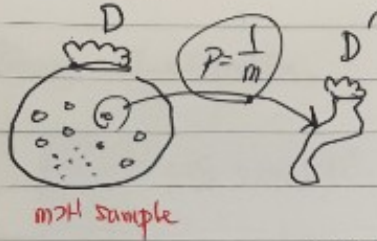


Train- Test



→ 결과 평균을!

③ Bootstrapping

한번 sample 1개씩 복사!

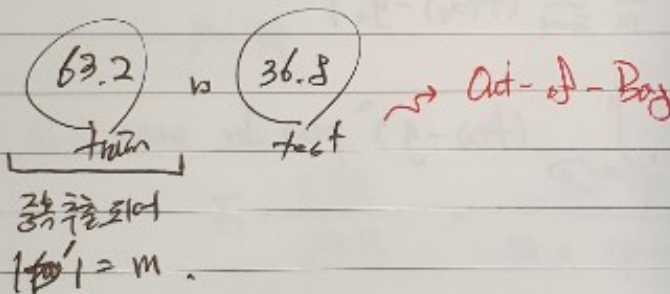
(복원 sample이 다시 복원될 수 없음!)

 $\times m$ 번 반복!

* 여러번 복원 sample로 만들고
원본으로 복원하지 않은 sample로 조합한 것!

$$\lim_{m \rightarrow \infty} \left(1 - \frac{1}{m}\right)^m = \frac{1}{e} \approx 0.368$$

\therefore 데이터가 충분히 클 때, D중 36.8%는 샘플로 D'로 들어가지 않음!
($m \rightarrow \infty$)



어제 쓰기에 좋을까? ① 데이터 세팅에 대해 잘룰 때

② 훈련/테스트 셋을 충분히 만들 때

+ 외에 다양한 D를 구성할 수 있는 매체에 (\therefore 학습자)

양상을 기법에 적용하기 좋음!!

But, 이 데이터 분포와 다른 가능성도

데이터가 충분히 많다면 Hold-out (CV)를 사용.

Date.

No.

④ 파라미터 튜닝과 학습 모델.

< 앙퍼리즘 파라미터, Hyper Parameter
모델 파라미터

→ 튜닝을 위해 **Validation Set** 을 사용하기로 함!

2.3 모델 성능 측정!

Performance Measure.

$$D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$$

f = learner. Compute $f(x)$, y

< Mean Squared Error > → **in Regression!** (시켜준)

$$E(f; D) = \frac{1}{m} \sum_{i=1}^m (f(x_i) - y_i)^2$$

Generally

$$E(f; D) = \int_{x \in D} (f(x) - y)^2 p(x) dx$$

< 정확도, 정합도 >

$$E(f; D) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(x_i) \neq y_i)$$

$$\text{acc}(f; D) = \frac{1}{m} \sum_{i=1}^m \mathbb{I}(f(x_i) = y_i) = 1 - E(f; D)$$

Generally

$$E(f; D) = \int_{x \in D} \mathbb{I}(f(x) \neq y) p(x) dx$$

$$\text{acc}(f; D) = \int_{x \in D} \mathbb{I}(f(x) = y) p(x) dx = 1 - E(f; D)$$

< Recall, Precision, F1-score >

In Binary-Classification,

	Actual	
	양성	음성
Predicted	양성	TP
	음성	FN
	양성	FP
	음성	TN

Confusion Matrix
확인하기 많이 !!

$$\text{Precision} = \frac{TP}{TP + FP}$$

양성은 (양성) 이라 할거임 맞는지

$$\text{Recall} = \frac{TP}{TP + FN}$$

Trade-off 관계 !!

⇒ 때문에, precision과 recall을 종합적으로 보아야 함

① BEP (Break-even point)

$$P = R$$

② F1-score

$$F_1 = \frac{2PR}{P+R} = \frac{2 \times TP}{m + TP - TN}$$

$$\frac{1}{F_1} = \frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right)$$

③ General F_β

$$F_\beta = \frac{(1+\beta^2) \times P \times R}{(\beta^2 \times P) + R}$$

$$\left(\frac{1}{F_\beta} = \frac{1}{1+\beta^2} \left(\frac{1}{P} + \frac{\beta^2}{R} \right) \right)$$

$\beta > 1 \Rightarrow R$ (재현율) 이 중요!

$\beta < 1 \Rightarrow P$ (정밀도) 이 중요!

Date.

No.

1) 다수의 이진 분류 Confusion Matrix를 얻거나,
다중분류의 결과?

① 모든 Confusion Matrix에 대해 P, R 계산 후 평균 계산.

$$\text{macro-P} = \frac{1}{n} \sum_{i=1}^n P_i = \text{mean}(P_i)$$

$$\text{macro-R} = \frac{1}{n} \sum_{i=1}^n R_i = \text{mean}(R_i)$$

$$\text{macro-F}_1 = \frac{2 \cdot \text{macro-P} \cdot \text{macro-R}}{\text{macro-P} + \text{macro-R}}$$

② 각 Confusion matrix의 원소값 (TP, TN, FP, FN) = 1 또는 0 계산.

$$\text{micro-P} = \frac{\overline{TP}}{\overline{TP} + \overline{FP}}$$

$$\text{micro-R} = \frac{\overline{TP}}{\overline{TP} + \overline{FN}}$$

* TP, TN, FP, FN 헷갈리지 말라!!

< ROC, AUC >

(cut point)

ROC (Receiver Operating Characteristics)

→ 적외선 레이저 신호 송수신기, 심리학과, 의학용 예언기 개발

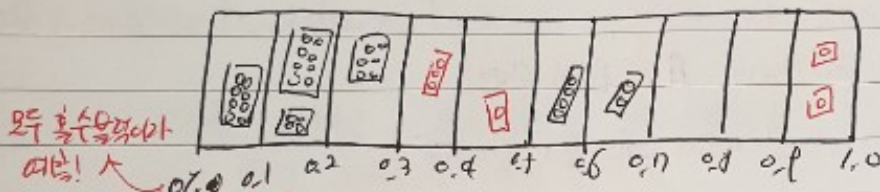
$$TPR \text{ (True-Positive rate)} = \frac{TP}{TP + FN}$$

$$FPR \text{ (False-Positive rate)} = \frac{FP}{TN + FP}$$

		predict	
		양성	음성
Actual	양성	TP	FN
	음성	FP	TN

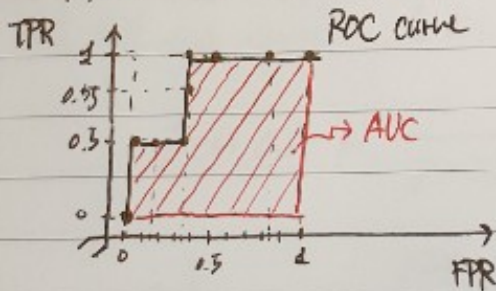
예제: 강태성 문제

양성 : P / 전체 : N



즉, 양성일지라도
예제!

양성 예측치	0%	10%	20%	30%	40%	50%	60%	70%	80%	90%	100%
TP 양성 수	4	4	4	4	3	2	2	2	2	2	0
FN 양성 수	0	1	3	4	4	4	5	6	6	6	6
FP 음성 수	0	0	0	0	1	2	2	2	2	2	4
FP 음성 수	6	5	3	2	2	2	1	0	0	0	0
TPR	1.0	1.0	1.0	1.0	0.75	0.5	0.5	0.5	0.5	0.5	0.0
FPR	1.0	0.34	0.5	0.34	0.34	0.34	0.16	0.0	0.0	0.0	0.0



※ 참고 !!

• 민감도 (Sensitivity)

= TPR

• 특이도 (Specificity)

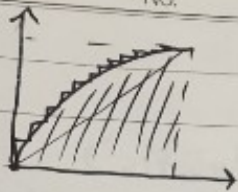
= 1 - FPR = TNR

(True Negative Rate)

Date.

No.

하버원!



$$AUC = \frac{1}{n} \sum_{i=1}^{n-1} (x_{i+1} - x_i) \cdot \frac{(y_i + y_{i+1})}{2}$$

$$\cancel{\frac{1}{n} \sum_{i=1}^{n-1} (x_{i+1} - x_i) y_{i+1}}$$

(y_i, y_{i+1}가 달라진 순간...? 있으니까?)

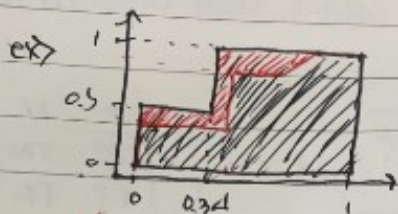
$$0.34 \times 0.3 + 0.68 \times 1 = 0.83$$

FPR이 변할 때
TPR은 변하는 경향이
있어.

↓
있을 수 있지!

threshold가 0.1 들어 들어 있으면

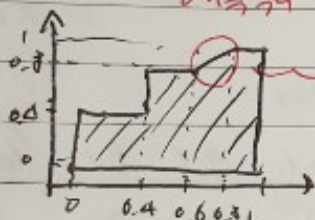
가 아닌 3번만 있을 수 있지!



이런 예제에서 0.1 ~ 0.2의

값을 threshold로 해서 바꾸면,

예제에서 낮아진다면 AUC값이 감소!



이 관 대응에 따라 같이 수식 사용!

4 AUC는 "생물 예측의 배열 순서 품질"을 고려.

m^+ 의 양성값, m^- 의 음성값, D^+, D^- ; 양/음성값의 집합

$$L_{rank} = \frac{1}{m^+ m^-} \sum_{x^+ \in D^+} \sum_{x^- \in D^-} \left(\mathbb{I}(f(x^+) < f(x^-)) + \frac{1}{2} \mathbb{I}(f(x^+) = f(x^-)) \right)$$

순서화

변경 후의

RDC 순서위의 변경, 거짓 양성률!

$$AUC = 1 - L_{rank}$$

< 비용 민감도 변화와 비용 곡선 >

In Real, '서로 다른 종류의 오차 가져오는 결과가 다른 현상'.

↙
아래한 서로 다른 종류의 cost에 대한 교점을 맞추기 위해
비균등 비용 unequal cost 사용.

• Binary Classification

cost matrix (비용 행렬)

$cost_{ij}$ = i 클래스 샘플이 j 클래스로 분류된 경우의 비용

		0 type	1 type	product
Real	0 type	0	$cost_{01}$	
	1 type	$cost_{10}$	0	

✓ 앞서 정의한 변화는 '오차 함수'로 직접 계산되고

오차 함수의 영향권에 대해 고려 X

✱

비용 민감도 cost-sensitive matrix

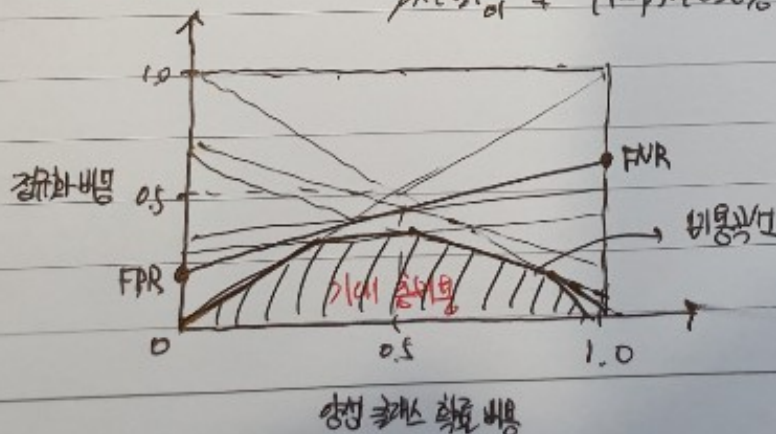
$$E(f; D; cost) = \frac{1}{m} \left(\sum_{x \in D} \mathbb{I}(f(x) \neq y_x) \times cost_{01} + \sum_{x \in D} \mathbb{I}(f(x) \neq y_x) \times cost_{10} \right)$$

✱

비용 곡선 cost curve (ROC 곡선 대체!)

$$P(+) cost = \frac{p \times cost_{01}}{p \times cost_{01} + (1-p) \times cost_{10}} \quad (가중치)$$

$$cost_{norm} = \frac{FNR \times p \times cost_{01} + FPR \times (1-p) \times cost_{10}}{p \times cost_{01} + (1-p) \times cost_{10}} \quad \text{where } FNR = 1 - TPR$$



Date.

No.

→ 학습 알고리즘 일반화 성능을 해석할 수 있는 이론적 도구.

2.5 Bias and Variance

편향-분산 문제 bias-variance decomposition.

$$\bar{f}(x) = E_D[f(x; D)] \quad ; \quad \text{학습 알고리즘의 기대 예측, 예측값}$$

$$\text{var}(x) = E_D[(f(x; D) - \bar{f}(x))^2] \quad \text{데이터셋 바뀔 때의 학습 성능 변화 측정}$$

$$\epsilon^* = E_D[(y_D - y)^2] \quad \text{기대 일반화 오차의 lower bound.}$$

$$\text{bias}^2(x) = E_D[(\bar{f}(x) - y)^2] \quad \text{학습 알고리즘의 적합 능력 평가.}$$

Assume that $E_D[y_D - y] = 0$.

$$\begin{aligned} E(f; D) &= E_D[(f(x; D) - y_D)^2] \\ &\stackrel{\text{일반화 기대 일반화 오차}}{=} E_D[(f(x; D) - \bar{f}(x) + \bar{f}(x) - y_D)^2] \\ &\stackrel{(\because \text{기대값})}{=} E_D[(f(x; D) - \bar{f}(x))^2] + E_D[(\bar{f}(x) - y_D)^2] \\ &\quad + E_D[2(f(x; D) - \bar{f}(x))(\bar{f}(x) - y_D)] \quad \begin{matrix} \because \bar{f}(x) = E_D[f(x; D)] \\ \text{데이터가 바뀔 때} \\ \text{일반화 오차} \end{matrix} \\ &= E_D[(f(x; D) - \bar{f}(x))^2] + E_D[(\bar{f}(x) + y - y - y_D)^2] \\ &= E_D[(f(x; D) - \bar{f}(x))^2] + E_D[(\bar{f}(x) - y)^2] \quad (\because \text{By Assumption}) \\ &\quad + E_D[(y - y_D)^2] + E_D[2(\bar{f}(x) - y)(y - y_D)] \\ &= \text{var}(x) + \text{bias}^2(x) + \epsilon^* \end{aligned}$$

Bias → bias-variance dilemma! (trade-off)