

Diagnostic Predictive Modeling

Date.

No. 2023. 07. 14 (Tue)

! 머신러닝을 하는 이유는? 자동이라서.
자동화된 process.

알고리즘은 너무 많아.

때때로 잘 되어있어서 쉽게 할 수 있어.

>

- 학습 dataset 훈련 (training)

- 학습 dataset 평가 (Validation)

- 학습 dataset 테스트 (test)

>

random_state = 42, 2014, 0...

stratify. ['survived']라서 0과 1이 test와 train에 고루
분포되게...

확률도 시뮬레이션. 신뢰성을 위해서.

> 알고리즘을 공부해야 하는 이유는 파라미터의 존재에 대해서 알아야 하기 때문.
파라미터가 어떻게 되느냐에 따라 성능이 좌우됨.

파라미터들을 관리하자.

> threshold를 0.5로 잡는 것이 아니라 조정해서 높은 성능을 더
올릴 수 있다.

KNN → K-Nearest Neighbours (KNN)

Neighbours 수가 바뀌면 정확도도 달라진다.

Random forest = 앙상블 ... 많은 모델이 합쳐짐.

여러개를 합쳐서 더 정확해진다.

앙상블이나 부스팅을 쓰고 2개를 더 합친다.. 더 정확도를 올린다.

이제는 다양한 모델이 많아야 한다.

그러나 cost 때문에 앙상블을 많이 안함.

Kaggle에서는 Score 때문에 많이 함.

모델을 몇 100개의 Stacking...

» 훈련 및 테스트 데이터가 변경되면 정확도도 변경된다.

정확도가 좋아지면 감소할 수 있다. 이를 극복하기 위해 교차검증을 사용한다.

christmas Arigato gojamas

Crossdresser
Yamanfunder

교차검증이 처음인 이야기가 바뀌어도 얼마나 유용한가?

Acc1	Acc2	Acc3	Acc4	Acc5
Model1	Model2	Model3	Model4	Model5
Learn	Learn	Learn	Learn	<u>predict</u>
Learn	Learn	Learn	<u>predict</u>	Learn
Learn	Learn	<u>predict</u>	Learn	Learn
Learn	<u>predict</u>	Learn	Learn	Learn
<u>predict</u>	Learn	Learn	Learn	Learn

$$(acc1 + acc2 + acc3 + acc4 + acc5) / 5$$

$$mean_acc = ?$$

$$Std_acc = ?$$

Y_c target 값 Survived.

⇒ Confusion Matrix

실제

0	491 TP	58 FN
1	95 FP	299 TN
	0	1

predict

→ 0은 사실
→ 1은 생존
→ 299는 생존 예외

실제 생존