

Supervised learning IV

모델 선택에 대한 문제. (모델 평가, 교차 검증, 부동 조정 선택, regularization)에 대해 소개.

여러 모델의 $f(x)$ 와 유사한 y 를 찾는 것
모든 모델에 대해 직접적으로 $f(x)$ 를 찾을 수
없기 때문에 훈련을 통해 $\hat{f} = \hat{f}(x)$ 를 구하게 됨
이 \hat{f} 가 진짜 모델의 $f(x)$ 와 유사한지
판별하기 위해서, Test set으로 검증할 것.
지도학습의 모델은 MSE를 줄이는 것을 목표로
함. 그러나, Train set의 MSE만을 고려하면
다른 데이터에는 적합하지 못한, Overfitting된
 $\hat{f}(x)$ 를 갖게 됨.

적성은 test set에 대해 최소 MSE를 갖는
모델 $\hat{f}(x)$ 를 찾아야 하는 것이고, 모델 적합성에는
train set만 사용해야 함.

모델 선택 기준 (선택화의 기준)

독립변수 d 에 따른 값으로 이 기준값이 최소화
혹은 최대화 되도록 학습시키는 기준을 제공한다

→ Mallows's C_p

$$C_p = \frac{1}{n} (RSS + 2d\sigma^2) \quad \rightarrow \quad \text{MSE}$$

→ Bayesian Information Criterion

$$BIC = \frac{1}{n} (RSS + \log n \cdot d\sigma^2) \quad \rightarrow \quad \text{MSE}$$

→ Adjusted R^2

$$Adj R^2 = 1 - \frac{RSS/(n-d-1)}{TSS/(n-1)} \quad \rightarrow \quad R^2$$

모델 선택 기준은 이론적 배경이 강하며,
전제 하에 관찰되지 않은 전체 집단의 MSE를
반영한다는 볼 수 있음.

→ 방정식으로 계산되며, 작은 값일수록
저산 가능함.

Cross validation (CV) : 교차 검증.

훈련세트 내에서 트레이닝 데이터를 샘플링함
모델 선택 기준이 비례 많은 데이터와 계산능력을
요구하며, 데이터에 대한 가절이 적음.
모세에는 교차검증으로 많이 사용함.

LOOCV (Leave One Out Cross Validation)

n 개의 샘플 중 1개를 제외한 상태로 학습하고,
제외한 1개의 샘플로 검증하는 방식.

K-fold CV

샘플을 k 개의 그룹으로 나누고 1개의 그룹을
제외한 상태로 학습 및 검증하는 방식.

LOOCV vs K-fold CV.

LOOCV가 항상 더 많이 적용함.

트레이닝 세트에 의존적이며, 더 유연함 (편향은
낮고, 분산은 높음)

K-fold는 저산의 횟수가 항상 더 적으며,
 k 를 n 으로 늘릴 경우, LOOCV와 동일함

변수 선택의 문제

전통적으로는 변수의 개수가 샘플에 비해 적어야
하는데는 오히려 변수 선택보다 많아지고 있다.
이런 상황에서는 전통적 학습법이 실패하기
때문에, 유익한 변수를 찾고, 고르는 작업이
필요하다.

→ 1개의 변수 중 1개이 변수만 선택
하는 것으로 모든 조합을 만들어보고
test set에 대한 최소의 MSE를 갖는
모델을 고르면 되지만, 현실적으로 어렵다.