

Overview

가계 함수로 많은 하위 범주를 포괄함

이 강좌에서는 가계 함수와 데이터 과학의 응용 영역

가계 함수로 사전에 정해진 특정 규칙을 따르는 것이 아니라, 데이터에 기반하여 의사결정하고 학습하는 머신을 만드는 것임.

가계는 함수로 표현 가능 $Y = f(X) + \epsilon$

Y 의 유형에 따라 함수의 형태가 달라짐.

$f(\cdot)$ 는 사전 정의되지 않고, X 에 따라 찾아져야 함
 X 는 입력 데이터. ϵ 는 풀 수 있는 노이즈임.

새로운 X 에 대한 Y 를 알지 못하면 $f(\cdot)$ 를 먼저 추정해야 함.

지도 학습에서는 특정 Y 를 통해 $f(\cdot)$ 를 모델링할 수 있음. $f(x)$ 로 알려진 Y 가 Y 와 차이가 적도록 모델링해야 함.

모차의 제곱 합을 최소화 하기 위한 식의 전개

$$E(Y - \hat{Y})^2 = E(Y - f(X))^2$$

$$= E[f(X) + \epsilon - f(X)]^2$$

$$= E[f(X) - f(X)]^2 + \text{Var}(\epsilon)$$

$\rightarrow E[\sim]$ 는 reducible error.

$\text{Var}(\epsilon)$ 는 irreducible error

\rightarrow 줄일 수 있는 에러를 줄여야 함!

회귀와 분류의 예측 MSE

$$\text{회귀 MSE} = \frac{1}{n} \sum (y_i - \hat{f}(x_i))^2$$

$$\text{분류 MSE} = \frac{1}{n} \sum I(y_i \neq \hat{f}(x_i))$$

Y 와 연관된 X 를 찾는 Inference

Effect size : X 증가 시 Y 의 증가를 증량화

Statistical significance : 통계 리스크로 유의성 증명

$f(\cdot)$ 를 추정 (estimate) 하는 방법.

모수적 (Parametric) : X 와 Y 의 특정 관계를 가정

\rightarrow 정해진 관계를 알 수 있으므로 선형성을 가짐.

예) 선형 회귀 = $\beta_0 + \beta_1 \text{나이} + \beta_2 \text{교육수준}$

비모수적 (Non-Parametric) : X 와 Y 사이에 관계가 없음을 가정

\rightarrow 주어진 유연성을 가진 모델을 만들어야 함.

유연성과 해석 가능성의 trade off 관계

유연성의 개념: 모델이 더 복잡하고, 통찰력이 있대야 함

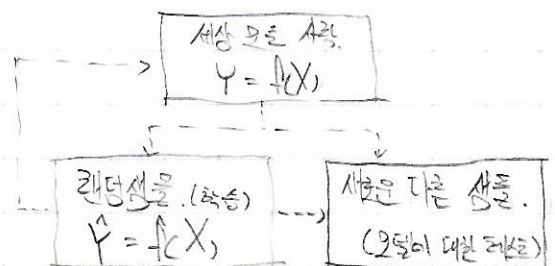
모델의 유연성이 높아질수록 해석 가능성을 떨어뜨림.

해석 가능성이 높으면, 복잡한 관계를 모델링할 수 없음.

비지도 학습 (자율 학습)은 Y 가 주어지지 않은 채로 $f(\cdot)$ 를 추정해야 함. 오히려 X 로 $f(\cdot)$ 를 추정함.

강화 학습은 동작에 따라 보상을 받고, 보상을 최대화하는 다음 조치를 찾는 것임. $Y = f(a|x)$

새로운 모델에 대한 특성을 찾는 것은 훈련에 사용된 X 가 아닌, 새로운 샘플을 통해 테스트되어야 함.



MSE는 train set과 test set에 대해 구할 수 있고, 최종적으로는 test set에 대한 MSE를 줄여야 함. 이 두 값을 통해 오버피팅, 언더피팅을 확인할 수 있음.