

Linear Regression 연속적인 데이터 $(-\infty, \infty)$

$$y = h_{\theta}(x) = w_0 + w_1 x_1 + w_2 x_2 + \dots$$

수학에서의 함수와 다르게 여기서는 w_0, w_1, w_2, \dots 이 마다 역할을 하고
데이터 (x, y) 를 넣어 가장 적합한 w_i 를 구하는 방식
손실 함수를 최소화

손실 함수

$$L(w) = \sum (\underbrace{y^{(i)}}_{\text{실제값}} - \underbrace{\hat{y}^{(i)}}_{\text{예측값}})^2$$

Logistic Regression 이분적인 데이터 (0 or 1)

Linear Regression의 치역과 동일해지려고 하면

$$\ln\left(\frac{p}{1-p}\right) = w_0 + w_1 x_1 + \dots$$

이후 나온 결과값 p 를 알아내기 쉽게 하면

$$p = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + \dots)}}$$

p : 1이 나올 확률 $1-p$: 0이 나올 확률.

손실 함수

$$L(w) = \sum (-y^{(i)} \log(\hat{y}^{(i)}) - (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$

학습 성능

K-fold 교차검증

S_1, S_2, S_3, S_4, S_5

S_1, S_2, S_3, S_4, S_5

S_1, S_2, S_3, S_4, S_5

S_1, S_2, S_3, S_4, S_5

S_1, S_2, S_3, S_4, S_5

데이터셋을 k 개로 나누어서 k 개는 테스트, 나머지는 훈련용으로 k 번 학습을 진행 후 정확도 비교

과적합

해결책

prepruning

기준점보다 낮은 향상을 보일 경우에는 더 진행하지 않는다.

postpruning

가지치는 것이 더 낮은 오류율을 이룰 경우, 가지 친다

Instance-Based Learning

본인 데이터와 비슷한 데이터들을 비교 후 결과 내린다