

# Supervised Learning I

다양한 회귀 분석법 (로컬 회귀, 다항식 회귀, 단계적 회귀, spline) 과 주성분 회귀, 부분최소제곱에 대한 내용

선형회귀의 단점 (다순하고 정확함)

비선형관계를 포착하기 어려움  $\rightarrow$  loess, spline

고차원 데이터에 적합하지 않음  $\rightarrow$  pcr, pls

Loess vs 일반 선형회귀

- Loess는 "에 비해 비모순적이다.
- "에 비해 유연성이 높다.
- "에 비해 많은 계산능력이 필요하다.
- "에 비해 많은 데이터가 필요하다.
- "에 비해 과대소조하지 않는다.

Loess (Local Regression) : 로컬 회귀

평평한 회귀를 local에서만 적용하는 개념



한 점을 기준으로 주위의 샘플로만 선형회귀를 단계적 회귀는 X의 전체범역을 구간으로 나누고, 수월하면, 한 점과 그 주변의 선형회귀선을 작구간의 값을 상수로 맞추는 아이디어임. 찾을 수 있음. 이를 모든 점에 대해 반복하면 곡선의 형태 (비선형)의 회귀선을 찾을 수 있음.

- 이것을 고려하는 것은 이상 가파를 수월 높은 가중치를 주는 개념임.

Loess의 RSS는 가중치를 계산 해야 함.

$$\rightarrow RSS_{loess} = \sum w_i (y_i - \beta_0 - \beta_1 x_i)^2$$

중심 가중치는 한점으로부터 근구간의 중심에 위치

$$\rightarrow w_i = \exp(-(x_i - x_0)^2 / 2\tau^2)$$

$\tau$  (tau)는 정규분포의 분산을 의미함.

tau가 크면 넓은 분포량이 되어, 정규분포에 멀리 떨어진 데이터포인트도 고려하게 됨.

tau가 작으면 좁은 분포량이 되어

"는 고려하지 않게 됨.

$\rightarrow$  loess의 유연성은  $\tau$ 에 의해 제어됨.

- 가중치 함수는 사각형의 형태를 띌 수 있음.

이렇게 되면 방위를 벗어난 데이터 포인트를 이에 고려하지 않게 됨.

Polynomial Regression : 다항 회귀

고차원 다항회귀는 선형회귀에 비해 유연성이 높음.

이는 회고차함에 의해 좌우된다. 이상치에

매우 민감하며, 오차간의 차이가 있다.

단계적 회귀는 X의 전체범역을 구간으로 나누고, 작구간의 값을 상수로 맞추는 아이디어임.

X를  $c_0, c_1, \dots, c_k$ 의 구간으로 분리함.

각 X에 대해서 구간의 값을 구함.

$$C_0(X) = I(X < c_0), C_1(X) = I(c_0 < X < c_1) \dots$$

I는 참 거짓 값으로 1, 0을 return 함.

위와 같은  $C_k$ 를 가중함수로 활용.

$$\hat{y}_i = \beta_0 + \beta_1 C_1(X_i) + \beta_2 C_2(X_i) + \dots + \beta_k C_k(X_i) + \epsilon_i$$

각 범위 내에서 계수인  $\beta$ 에 따라  $y_i$ 가 결정됨



절단점의 갯수에 따라 유연성이 결정되며,

통상 3~5개를 사용함.

$\rightarrow$  3개를 사용할 경우 사분위수와 같이

Cut point를 설정함.