

- 회귀분석에 사용되는 데이터는 그 자체로 사용되기 보다는 스케일이나 함수변환 등의 처리 과정을 거쳐야 경우가 많다. 왜냐하면 과정은 공정한 함수의 도출을 향상시키기 위해 데이터 간의 관계를 선형화하기 위해 바꾸기 위해서 사용된다.

- 조건수 (Condition number)

- 조건수는 공변행렬의 $X^T X$ 의 가장 큰 고유값과 가장 작은 고유값의 비율을 뜻한다.

$$\text{Condition number} = \frac{\lambda_{\max}}{\lambda_{\min}}$$

- 조건수가 크면 예측을 계산할 때 오차가 미치는 영향이 커진다.

예) $Ax = b$

조건수 작은 경우

- *행렬 A가 단위행렬이면 조건수는 가장 작은 값으로 조건수가 1이다.

$$\text{Cond}(I) = 1$$

조건수 큰 경우

- *조건수가 크면 오차의 왜곡 및 오차가 전체 다른 값을 가진다. 따라서 조건수가 크면 회귀분석은 사용은 예측값의 왜곡을 귀하게 된다.

- 회귀분석과 조건수

- 회귀분석에 조건수가 커지면 예측은 크게 두가지가 있다.

1. 변수들의 단위 차이로 인해 변수의 스케일이 크게 달라지는 경우. 이 경우에는 스케일링(Scaling)으로 해결한다.

2. 다중공선성 즉, 상관관계가 큰 독립변수들이 있는 경우, 이 경우에는 변수 선택이나 PCA를 사용한 차원축소로 해결한다.

스케일링
독립변수의 단위가 제각각이므로 scale 함수를 통해 부채를 비슷하게 맞추어 줄 수 있다.
Scale 함수는 표준편차로 표준화한다.

변수선택
독립변수나 종속변수가 상당히 높은 차원을 포함하는 경우 / 독립변수 간의 상관관계가 높거나 종속변수가 비선형 관계를 보이는 경우
관계를 가장 좋은 방법으로 모델링 / 종속변수의 예측이 비선형 관계를 보이는 경우
조건수가 높을 경우 이러한 회귀분석이 향상될 수 있다.