

Categorical predictors

예를 들어 $x \in \{0, 1\}$ 인 경우엔 쉽게 처리가 가능하다.

$y = \beta_0 + \beta_1 x$ 을 라면 β_1 은 $x=1$ 이면 나타내기

$$\begin{pmatrix} y = \beta_0 \\ y = \beta_0 + \beta_1 \end{pmatrix} \text{ 두 식을 비교 가능하도.}$$

하지만 $x \in \{0, 1, 2\}$ 인 경우 1과 2를 비교해야한다. 그러므로

중에서 1과 2의 대조를 비교하면 안된다.

⇒ 결국 0, 1 또한 표현해야한다.

A, B, C가 존재하는 경우

$$x_1 = \begin{cases} A & \text{if } x=1 \\ \text{not } A & \end{cases}$$

$$x_2 = \begin{cases} B & \text{if } x_2=1 \\ \text{not } B & \text{else} \end{cases}$$

그래서 (x_1, x_2) 을 사용해보면

$$\begin{pmatrix} (1, 0) \Rightarrow A \\ (0, 1) \Rightarrow B \\ (0, 0) \Rightarrow C \end{pmatrix} \text{ 을 만들 수 있다.}$$

이제 $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2$ 로 표현이 가능해진다.

A: $y = \beta_1 + \beta_1 \Rightarrow C$ 대비 사용량 # 120 $\Rightarrow C$ 보다 A가 120더 사용

B: $y = \beta_0 + \beta_2 \rightarrow C$ 대비 사용량 # -3 $\Rightarrow C$ 보다 B가 -3더 사용

C: $y = \beta_0 \Rightarrow$ 사용량. # 500

아예한 방식은 category가 k개인 경우 (k-1)개의 변수를 도입한다.

↳ 여기서는 Bias가 캐리리(이 없다).

< one-hot-encoding >

모든 category 마다 변수를 할당한다. (여기서는 Bias가 특정 캐리리가 X)

Modeling Interaction

$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$ 처럼 Modeling 하면 X_1 과 X_2 를 독립적으로 본다.

But X_1 과 X_2 를 동시에 고려하면 상호작용을 알 수 있다.

(synergy, interaction effect)

$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \epsilon$ 로 modeling이 가능하다.

Model 공식 자체는 Linear Regression과 유사.

1) X_1 과 X_2 의 상충관계가 있으면 자연스럽게 β_3 는 0으로 나올 것이다.

"Hierarchy for interaction"

⇒ 우리가 interaction term을 넣기 위해서는 반드시 단일 term을 먼저 넣어줘야 한다. 예를 들어) $Y = \beta_0 + \beta_1 X_1 X_2$ (X)

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 (0)$$

반드시 modeling의 완전성이 깨지고, 해석이 어려워진다.

추가로 $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1^2 + \beta_4 X_2^2$ 도 가능하다. 이는 선형 모델의 구조를 유지하면서 비선형성을 추가해주는 역할을 한다.



Feature selection

ML에서 β 를 구하는 것 $\hat{\beta} = (X^T X)^{-1} X^T y$ 를 통해서 구하게 되는데, $(X^T X)^{-1}$ 이 존재하지 않거나 느릴 수 있다. \Rightarrow 특정 feature 제거가 필요.

그러고 영향력 ↓인 β_n 이 존재할 수도 있다.

pre-select

< 1. Best subset selection > $O(2^p)$

n 개의 feature 중에서 k 개를 select 하는 경우 가능한 nC_k 개의 조합을 모두 고려.

ex) step1

step2

step3

--->

$$y = 0.8$$

$$y = 0.2 + 0.1x_1$$

$$y = 0.1 + 0.4x_0 + 0.1x_2$$

* curse of dimensionality \Rightarrow 차원 ↑일수록 연산량 ↑이지만 gain ↓된다.

step-select

< 2. Forward stepwise > $O(p^2)$ But, best solution이 아닐 수 있다.

이전 step이 select 된 것들만 추가. (계수는 change 가능)

step1.

step2

step3

→ -----

$$y = 0.8$$

$$y = 0.8 + 0.4x_1$$

$$y = 0.8 + 0.3x_1 + 0.1x_2$$

< 3. Backward stepwise > $O(p^2)$

모든 N 개에서 1개씩 빼야한다 $\Rightarrow (N-1) \rightarrow (N-2) \rightarrow \dots \rightarrow 1$ 까지 가라.

단, $N > p$ 는 만족해야한다.

< 4. Stage wise >

2와 동일. But coefficient까지도 fix시킨다. But 1번은 변수를 라식시킬 수 있다 그래서 이를 통해서 coefficient를 조절 가능해진다.

Choosing the optimal model

stepwise selection에서 변수의 수를 ↑갓 수록 성능이제는 올라간다.

But 우리는 data의 "일반성" 특징/패턴을 파악하기 위한 것. → validation data로 검증.

cross validation (train : valid : test).

k-fold cross validation.