

Practice-HW

DM Process / Variable Selection

논문참조

강세정, 김규현, 신재욱, 장현우



01

Research Questions

Introduction to RQ
Data

02

Research Question 1

EDA & VIF
Regression (OLS)
Variable Selection
Shrinkage Comparison & Result

03

Research Question 2

Chi-square test
Result

Research Questions

Contents

DM Process /
Variable Selection

01. Introduction to RQ

- (1) RQ1 : 시민들의 여가생활과 문화생활에 미치는 날씨 요인은 무엇일까? → 따릉이, 경복궁 이용객 수의 수요 예측을 시도해보고자 함
- (2) RQ2 : 강수 여부가 대기 환경에 어떤 영향을 미칠까? → 강수 여부와 미세먼지, 초미세먼지의 독립 여부를 확인해보고자 함

02. Data

(1) RQ1

X(Numerical Variable)

평균기온
(°C)

최고기온
(°C)

일강수량
(mm)

최대풍속
(m/s)

평균풍속
(m/s)

최소상대
습도(%)

평균상대
습도(%)

가조시간
(hr)

합계일조
시간(hr)

1시간 최다
일사량
(MJ/m2)

합계 일사량
(MJ/m2)

일 최심
적설(cm)

평균전운량
(1/10)

9-9강수
(mm)

안개 계속
시간(hr)

X(Categorical Variable)

강수여부
(0, 1)

적설여부
(0, 1)

안개여부
(0, 1)

휴일여부
(0, 1)

미세먼지 기준
(1, 2, 3, 4)

초미세먼지 기준
(1, 2, 3, 4)

Y1

일일 서울시 공공자전거
(따릉이) 대여 건수

Y2

일일 경복궁 관람객 수

「대기환경보전법」에 따른 미세먼지, 초미세먼지 기준

물질	단위	산정기준	등급			
			총유(1)	보통(2)	나쁨(3)	매우나쁨(4)
미세먼지	μg/m³	24시간	0~30	31~80	81~150	151~
초미세먼지	μg/m³	24시간	0~15	16~35	36~75	76~

- **p-value** 0.05 이상인 변수들을 제외하면 통계적으로 유의미하다고 판단

Research Question 1

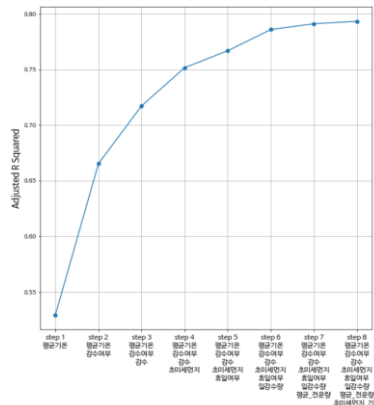
Contents

DM Process / Variable Selection

03. Variable Selection

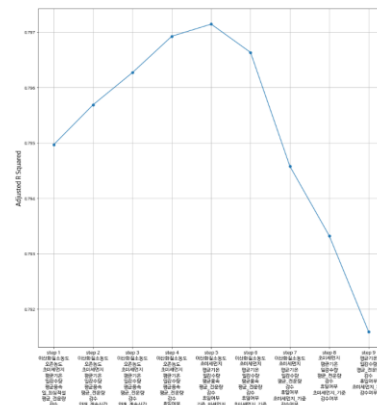
(1) Y1 : 일일 서울시 공공자전거(따릉이) 대여 건수

① 전진 선택법(Forward Selection)



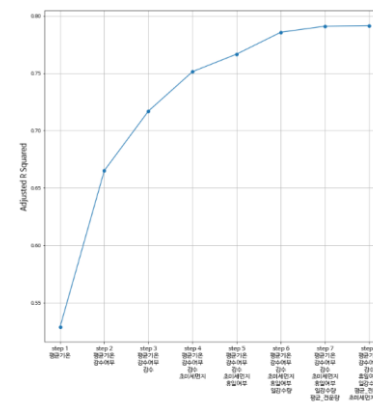
- 총 8개의 변수가 선택됨
- **F-statistics** = 140.6, **prob(F-statistics)** = 6.35e-94로 도출된 회귀식이 회귀분석 모델 전체에 대해 통계적으로 유의
- **Adj. R-squared** = 0.799으로, 전체의 80% 정도를 설명하고 있음
- **p-value** 0.05 이상인 변수들을 제외하면 통계적으로 유의미하다고 판단

② 후진 선택법(Backward Elimination)



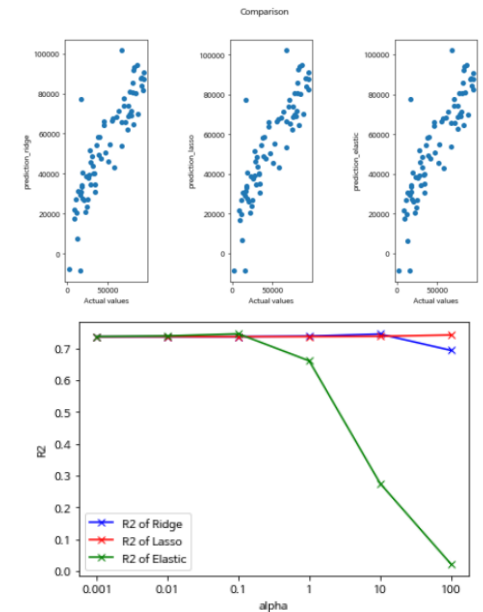
- 총 7개의 변수가 선택됨
- **F-statistics** = 158.9, **prob(F-statistics)** = 2.58e-94로 도출된 회귀식이 회귀분석 모델 전체에 대해 통계적으로 유의
- **Adj. R-squared** = 0.797으로, 전체의 80% 정도를 설명하고 있음
- **p-value** 0.05 이상인 변수들을 제외하면 통계적으로 유의미하다고 판단

③ 단계별 선택법(Stepwise Selection)



- 총 7개의 변수가 선택됨
- **F-statistics** = 158.9, **prob(F-statistics)** = 2.58e-94로 도출된 회귀식이 회귀분석 모델 전체에 대해 통계적으로 유의
- **Adj. R-squared** = 0.797으로, 전체의 80% 정도를 설명하고 있음
- **p-value** 0.05 이상인 변수들을 제외하면 통계적으로 유의미하다고 판단

④ Ridge & Lasso & Elastic Net



- 큰 alpha 값에 대해서는 Lasso가 가장 우수
- 모델별 최적 파라미터 상에서의 성능 비교
- R-squared** : Linear = 0.74, Ridge = 0.74, Lasso = 0.74, Elastic Net = 0.74
- MSE** : Linear = 202736114.02, Ridge = 198467211.98, Lasso = 198764427.26, Elastic Net = 199819289.10
- MAE** : Linear = 11019.17, Ridge = 10962.76, Lasso = 10940.26, Elastic Net = 10978.17

Research Question 1

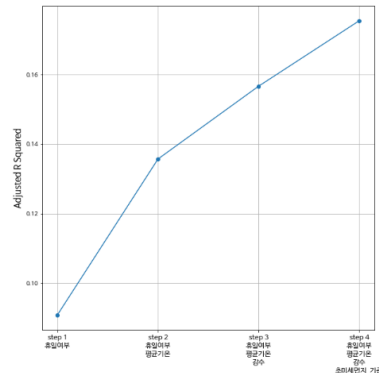
Contents

DM Process / Variable Selection

03. Variable Selection

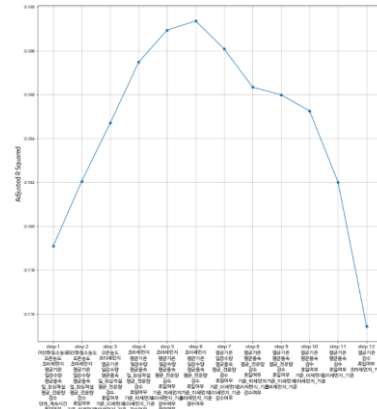
(2) Y2 : 일일 경북공 관람객 수

① 전진 선택법(Forward Selection)



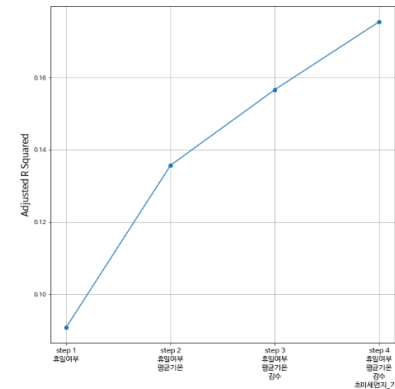
- 총 4개의 변수가 선택됨
- **F-statistics** = 16.48, **prob(F-statistics)** = 3.63e-12로 도출된 회귀식이 회귀분석 모델 전체에 대해 통계적으로 유의
- **Adj. R-squared** = 0.187으로, 전체의 19% 정도를 설명하고 있음
- **p-value** 0.05 이상인 변수들을 제외하면 통계적으로 유의미하다고 판단

② 후진 선택법(Backward Elimination)



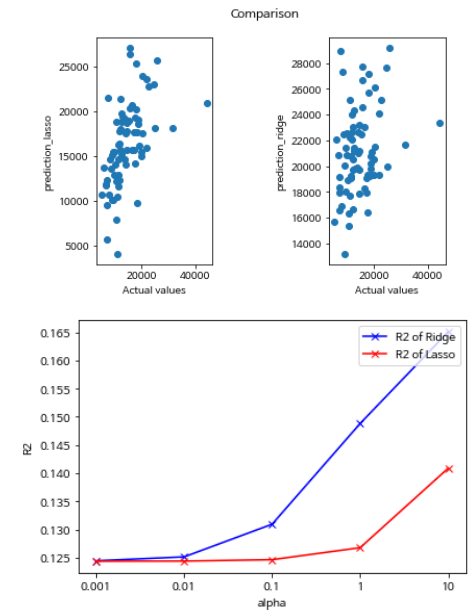
- 총 4개의 변수가 선택됨
- **F-statistics** = 16.48, **prob(F-statistics)** = 3.63e-12로 도출된 회귀식이 회귀분석 모델 전체에 대해 통계적으로 유의
- **Adj. R-squared** = 0.187으로, 전체의 19% 정도를 설명하고 있음
- **p-value** 0.05 이상인 변수들을 제외하면 통계적으로 유의미하다고 판단

③ 단계별 선택법(Stepwise Selection)



- 총 4개의 변수가 선택됨
- **F-statistics** = 16.48, **prob(F-statistics)** = 3.63e-12로 도출된 회귀식이 회귀분석 모델 전체에 대해 통계적으로 유의
- **Adj. R-squared** = 0.187으로, 전체의 19% 정도를 설명하고 있음
- **p-value** 0.05 이상인 변수들을 제외하면 통계적으로 유의미하다고 판단

④ Ridge & Lasso Regression



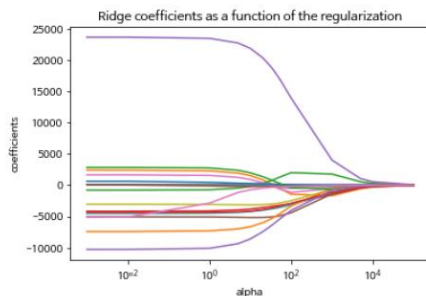
- Ridge의 성능이 Lasso보다 우수
- train accuracy : 0.22 / test accuracy : 0.12
- alpha 값의 증가에 따라 성능의 차이가 커짐

Research Question 1 / 2

Contents

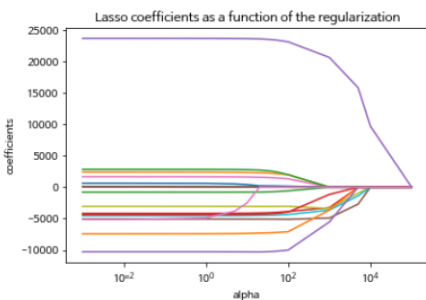
DM Process /
Variable Selection

04. Shrinkage Comparison & Result



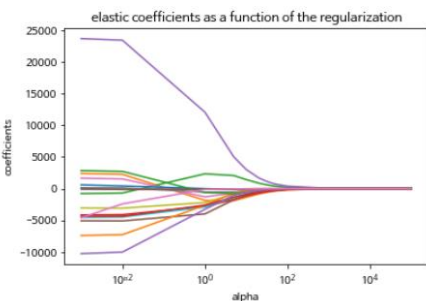
- Y1(따릉이)에 대한 회귀 모델에서 R-squared가 약 0.8 정도로 유의미하게 나와 예측 성능에 있어 우수하다고 판단함

- 다양한 변수 선택법과 더불어 Ridge, Lasso, Elastic Net에 대해 성능을 확인하였고, Shrinkage 효과를 비교 확인함



- Y1(따릉이) 예측에 있어 초기에 우리가 직접 수집한 변수들 중 기온, 강수, 초미세먼지 등의 변수들이 유의미함을 확인함

- Y2(경복궁) 예측 성능은 Y1(따릉이)에 비해 낮지만, 선택된 변수들은 동일함을 확인할 수 있었음



- 분석 이전에는 문화 활동과 여가 활동에 유사한 날씨 변수가 영향을 미칠 것으로 생각하였으나, Y2(경복궁)의 예측력은 Y1(따릉이)에 비해 매우 낮아 여가 활동에 대해서는 날씨 변수 이외에 추가적인 다른 요인을 고려해야 할 필요가 있음

Research Question 2

01. Chi-square Test

(1) 강수여부 - 미세먼지 독립 여부 판단

- **chi-square statistic** = 26.906, **p-value** = 0.000006 < 0.05

(2) 강수여부 - 초미세먼지 독립 여부 판단

- **chi-square statistic** = 14.599, **p-value** = 0.002194 < 0.05

02. Result

H_0 (귀무 가설) = 강우 와 (초)미세먼지는 서로 독립이다

H_A (대립 가설) = 강우 와 (초)미세먼지는 서로 독립이 아니다

→ p-value가 기준치 (0.05)보다 작아 귀무 가설을 기각하고 대립 가설 채택

Chi-square 검정 중 독립성 검정 결과, 상식과 동일한 결론을 얻을 수 있었으며, 자료가 다항분포나 이항분포를 따른다고 가정하며, 기대도수가 5보다 커야 한다는 Chi-square 검정의 가정 만족

Data Source

Contents

DM Process /
Variable Selection

[데이터 출처]

(1) 서울 열린 데이터 광장

- 서울시 일별 평균 대기오염도 정보 <http://data.seoul.go.kr/dataList/OA-2218/S/1/datasetView.do#>
- 서울시 공공자전거 이용현황 <https://data.seoul.go.kr/dataList/OA-14994/F/1/datasetView.do>

(2) 공공데이터포털

- 문화재청 궁능유적본부_4대궁 관람객 수 현황 <https://www.data.go.kr/data/15040885/fileData.do?recommendDataYn=Y>

(3) 기상청

- 일별 종관기상관측(ASOS) 자료 <https://data.kma.go.kr/data/grnd/selectAsosRltmList.do?pgmNo=36>

(4) 2016 ~ 19년 공휴일 데이터

- <https://superkts.com/day/holiday/2019>

[코드 출처]

(1) VIF 계산 코드

- <https://signature95.tistory.com/m/18>