

Practice-HW

Logistic Regression / Model Assessment

논문참조

강세정, 김규현, 신재욱, 장현우



01

Introduction

Research Objective
Data

02

Body

Preprocessing & EDA
Logistic Regression
Evaluation
Result

03

Conclusion

Result
Limitaion
Development

Introduction

Contents

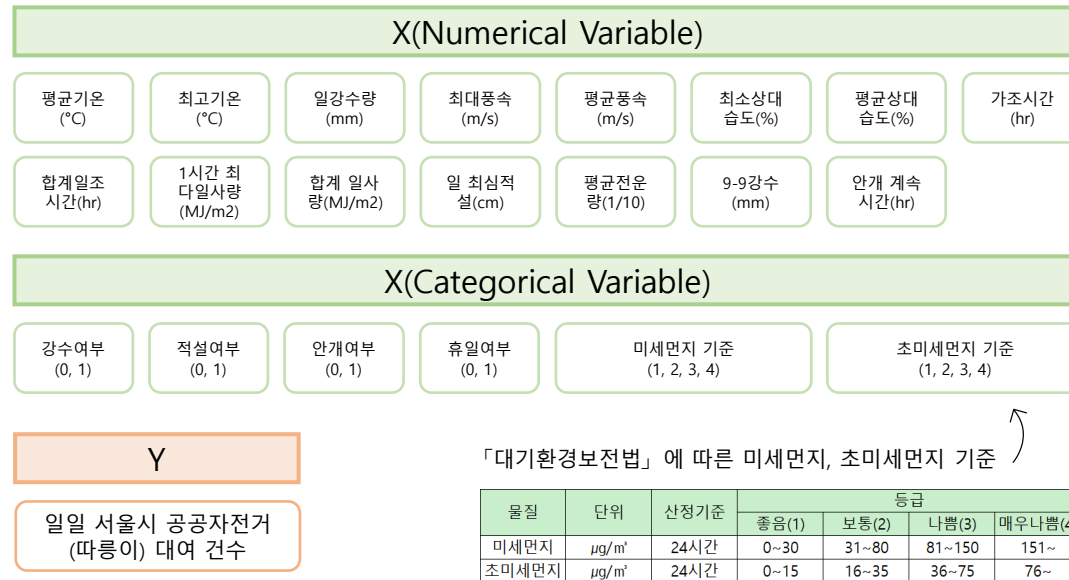
Logistic Regression
/ Model Assessment

01. Research Question

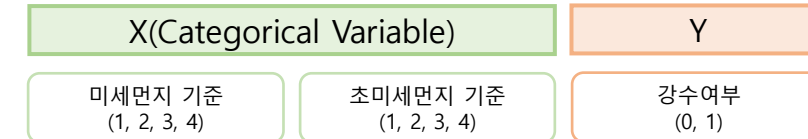
날씨 요인에 따라 따릉이 이용객 수가 평소보다 많은 지 여부를 예측해보자

02. Data

(1) RQ1



(2) RQ2



(3) 추가 설명

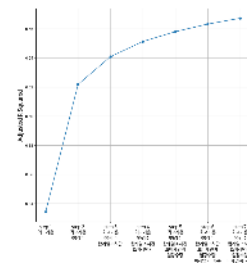
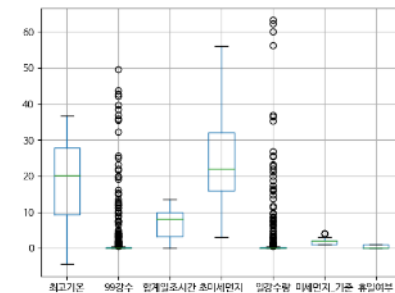
데이터는 서울 열린 데이터 광장, 공공데이터포털, 기상청 등에서 직접 수집함
데이터 기간 : 2019.01.01 – 2019.12.31

Y 값은 대여 건수 평균값 이상 / 이하로 Binary 데이터로 변환함
→ Binary Logistic Regression 수행 (해석 가능)

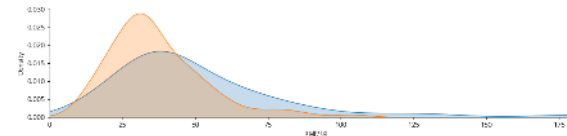
	여객차 승하차 수합계	모노노노	여객차 승하차 수합계	대형버스	중형버스	경전철 승하차 수합계	회고귀환	알뜰교통 카드 이용수합	비대면출 발... 수합계	경전철 승하차 수합계	99구분	간접적기 타사건	대형버스 승하차 수합계	중형버스 승하차 수합계	공수배부	학원배부	단체배부	유동배부	미등록	
count	365.000000	365.000000	365.000000	365.000000	365.000000	365.000000	365.000000	365.000000	365.000000	-	365.000000	365.000000	365.000000	365.000000	365.000000	365.000000	365.000000	365.000000	365.000000	
mean	0.003026	0.024042	0.070411	0.003345	432.6148	26.326787	13.959894	18620.192	2.241918	4.266049	-	4.899178	2.241918	0.023233	17.78082	2.002740	0.268493	0.109599	0.317808	52346.054979
std	0.013131	0.012714	0.171775	0.000859	24.099918	18.804996	10.164196	10335.448	1.815206	1.072120	-	2.979421	7.505520	0.295512	0.586607	0.716749	0.443784	0.104252	0.466264	27256.890666
min	0.000000	0.003000	0.200000	0.002000	8.000000	3.000000	7.900000	4.550000	0.000000	1.800000	-	0.000000	0.000000	0.000000	1.000000	1.000000	0.000000	0.000000	0.000000	26046.000000
25%	0.019750	0.014000	0.400000	0.003000	28.000000	16.000000	4.300000	9.300000	0.000000	3.500000	-	2.300000	0.000000	0.000000	1.000000	2.000000	0.000000	0.000000	0.000000	26046.000000
50%	0.027000	0.023000	0.400000	0.003000	32.000000	22.000000	14.700000	20.000000	0.000000	4.200000	-	4.300000	0.000000	0.000000	2.000000	2.000000	0.000000	0.000000	0.000000	54142.000000
75%	0.039000	0.048000	0.600000	0.004000	51.000000	32.000000	23.000000	27.800000	0.200000	4.900000	-	7.500000	0.200000	0.000000	2.000000	2.000000	1.000000	0.000000	0.000000	78874.000000
max	0.071000	0.059000	1.300000	0.006000	183.000000	117.000000	31.600000	36.800000	63.200000	11.800000	-	10.000000	49.600000	53.300000	4.000000	4.000000	1.000000	1.000000	1.000000	95959.000000

8 rows × 20 columns

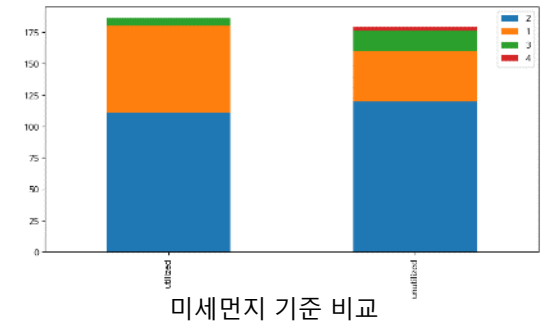
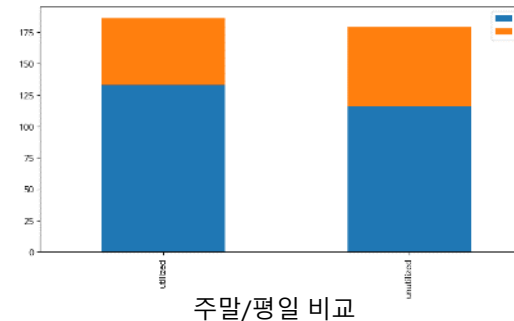
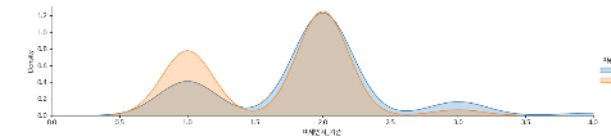
VIF는 모두 3을 크게 초과하지 않음 : 다중공선성 X



	variables	VIF
0	최고기온	3.608792
1	99강수	2.099618
2	합계일조시간	3.671065
3	초미세먼지	1.987403
4	일강수량	2.030607



추가로, 범주형 변수별 빈도 비교 결과 주말보다 평일에 따름이 대여가 더 활발하며, 미세먼지는 1(매우 좋음)일 때 높은 대여, 3,4(나쁨, 매우 나쁨)일 때 낮은 대여 빈도가 확인됨. 2(중음)일 때에는 별다른 차이 없음



두 가정은 추후에 확인

02. Logistic Regression

Undersampling을 통해 비율 맞춤	0	112
	1	112

```
Optimization terminated successfully.  
Current function value: 0.232170  
Iterations 33
```

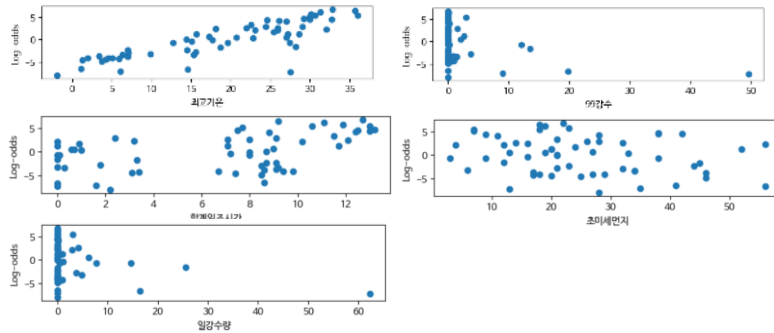
Body

Contents

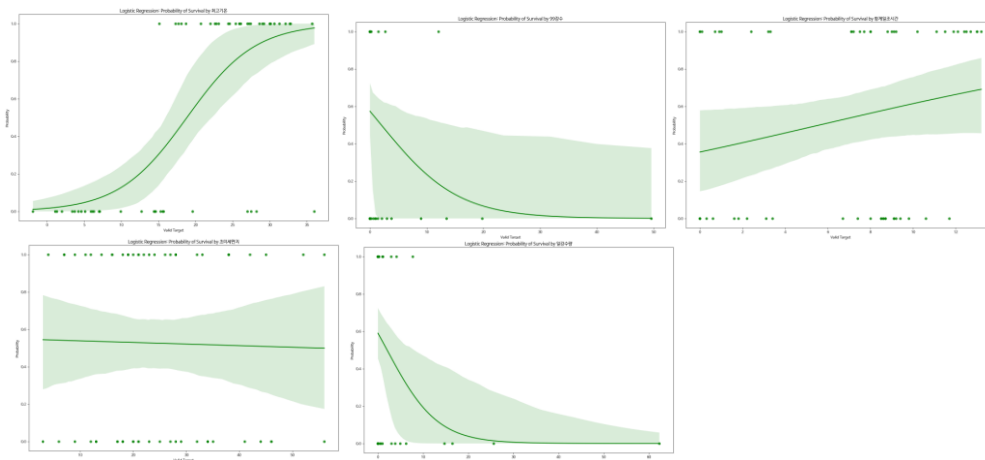
Logistic Regression / Model Assessment

03. Evaluation

(1) 가정 확인: 3. 독립 변수와 log odds 간의 linearity



99강수, 초미세먼지, 일강수량 변수는 비선형 패턴을 보여주므로 로짓 선형성 가정이 위반됨



(2) 모델 평가

최적의 cut-off value 선정

0.1 ~ 0.9 비교 결과, 0.2 가 최적

```
0.1
[1] 71
[0] 11
Accuracy: 0.8813559322033698
Specificity: 0.95
Precision: 0.91549614642450
Recall: 0.9
F1 score: 0.90875027467261

0.2
[1] 61
[0] 11
Accuracy: 0.8813559322033698
Specificity: 0.95
Precision: 0.91549614642450
Recall: 0.9
F1 score: 0.90875027467261

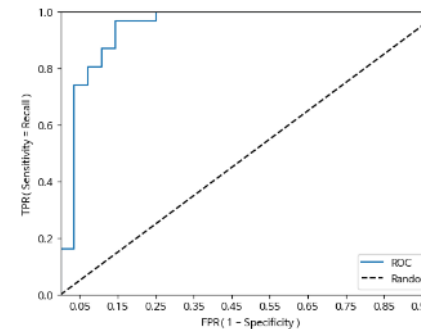
0.3
[1] 61
[0] 11
Accuracy: 0.8813559322033698
Specificity: 0.95
Precision: 0.91549614642450
Recall: 0.9
F1 score: 0.90875027467261

0.4
[1] 61
[0] 11
Accuracy: 0.8813559322033698
Specificity: 0.95
Precision: 0.91549614642450
Recall: 0.9
F1 score: 0.90875027467261
```

- 해석

- 22개의 샘플이 실제 Negative 클래스이고 예측도 Negative로 정확하게 예측
- 6개의 샘플이 실제 Negative 클래스이지만 예측은 Positive로 잘못 예측
- 1개의 샘플이 실제 Positive 클래스이지만 예측은 Negative로 잘못 예측
- 30개의 샘플이 실제 Positive 클래스이고 예측도 Positive로 정확하게 예측

ROC Curve 해석

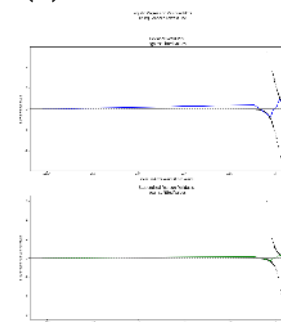


이 모델의 정확도는 88.14%로, 대부분의 샘플을 정확하게 예측하고 있음
특이도는 78.57%로, 실제 Negative 클래스를 정확하게 예측한 비율이 높음
정밀도는 83.33%로, Positive로 예측한 샘플 중에서 실제 Positive 클래스인 비율이 높음
재현율은 96.77%로, 실제 Positive 클래스 중에서 정확하게 예측한 비율이 높음
F1 점수는 89.55%로, 정밀도와 재현율의 조화 평균이 높음

종합적으로, 전반적으로 높은 성능을 보임

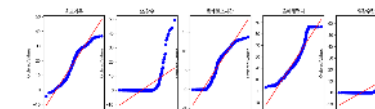
TPR (Sensitivity = Recall) 과 FPR (1- Specificity) 의 관계가
왼쪽과 같이 좌측 상단에 밀집한 그래프로 나옴
결과가 잘 나옴

(3) 가정 확인: 6. 잔차간 등분산성이 존재해야 함



실선은 잔차의 추세를 나타내며, 실선이 y = 0 의 intercept를 가진 상태로 수평으로 근사되는 것으로 보아, 모델의 불완정성이 낮음
그러나, 잔차가 잔차가 무작위로 배치되지 않았으므로, 잔차의 등분산성이 만족하지 않음

추가로, qqplot과 bartlett 검정으로도 확인



bartlett 검정 결과,
p-value < 0.05 → 귀무가설 기각

Results

Contents

Logistic Regression /
Model Assessment

Model.summary

Logit Regression Results						
Dep. Variable:	따릉이	No. Observations: 224				
Model:	Logit	Df Residuals: 214				
Method:	MLE	Df Model: 9				
Date:	Wed, 19 Apr 2023	Pseudo R-squ.: 0.6650				
Time:	13:48:40	Log-Likelihood: -52.006				
converged:	True	LL-Null: -155.26				
Covariance Type:	nonrobust	LLR p-value: 1.423e-39				
	coef	std err	z	P> z	[0.025	0.975]
const	-5.6014	1.238	-4.523	0.000	-8.029	-3.174
최고기온	0.3148	0.044	7.145	0.000	0.228	0.401
99강수	-0.1144	0.038	-3.030	0.002	-0.188	-0.040
합계일조시간	0.2207	0.079	2.800	0.005	0.066	0.375
초미세먼지	-0.0644	0.030	-2.124	0.034	-0.124	-0.005
일강수량	-0.0607	0.056	-1.090	0.276	-0.170	0.048
미세먼지_기준_2	0.6322	0.758	0.834	0.404	-0.854	2.118
미세먼지_기준_3	1.2346	1.473	0.838	0.402	-1.652	4.121
미세먼지_기준_4	-94.5226	9.32e+20	-1.01e-19	1.000	-1.83e+21	1.83e+21
휴일여부_1	-1.0285	0.581	-1.770	0.077	-2.168	0.111

Model.summary2

Model:	Logit	Pseudo R-squared: 0.665				
Dependent Variable:	따릉이	AIC: 124.0123				
Date:	2023-04-19 13:48	BIC: 158.1288				
No. Observations:	224	Log-Likelihood: -52.006				
Df Model:	9	LL-Null: -155.26				
Df Residuals:	214	LLR p-value: 1.4230e-39				
Converged:	1.0000	Scale: 1.0000				
No. Iterations:	33.0000					
	Coef.	Std.Err.	z	P> z	[0.025	0.975]
const	-5.6014	1.2384	-4.5232	0.0000	-8.0285	-3.1742
최고기온	0.3148	0.0441	7.1453	0.0000	0.2285	0.4012
99강수	-0.1144	0.0378	-3.0297	0.0024	-0.1884	-0.0404
합계일조시간	0.2207	0.0788	2.7995	0.0051	0.0662	0.3752
초미세먼지	-0.0644	0.0303	-2.1239	0.0337	-0.1239	-0.0050
일강수량	-0.0607	0.0557	-1.0897	0.2758	-0.1698	0.0485
미세먼지_기준_2	0.6322	0.7582	0.8338	0.4044	-0.8539	2.1183
미세먼지_기준_3	1.2346	1.4729	0.8382	0.4019	-1.6522	4.1214
미세먼지_기준_4	-94.5226	931596430646428434432.0000	-0.0000	1.0000	-1825895452193066057728.0000	1825895452193066057728.0000
휴일여부_1	-1.0285	0.5812	-1.7695	0.0768	-2.1676	0.1107

Odds ratio(승산비)

```
const          3.692859e-03
최고기온       1.369998e+00
99강수         8.918823e-01
합계일조시간   1.246932e+00
초미세먼지     9.375850e-01
일강수량       9.411214e-01
미세먼지_기준_2 1.881733e+00
미세먼지_기준_3 3.437082e+00
미세먼지_기준_4 8.899113e-42
휴일여부_1     3.575596e-01
dtype: float64
```

미세먼지_기준_2 : 미세먼지_기준_1보다
odds가 1.88배 증가
미세먼지_기준_3 : 미세먼지_기준_1보다
odds가 3.44배 증가

휴일여부_1 : 휴일여부_0보다 odds가
0.357배 증가 (64.3% 감소)

Null deviance : $P > |z|$: 0.05 미만이면 기각, x가 y를 설명하는 변수임을 의미

최고기온 / 99강수 / 합계일조시간 / 초미세먼지 : 귀무가설 기각, 따릉이 이용객 수를 설명하는 유의미한 변수

일강수량, 미세먼지_기준 (2,3,4), 휴일여부 : 귀무가설 기각할 수 없음, 따릉이 이용객 수를 설명한다고 할 수 없음

Logit(따릉이 = 1) = $-5.6014 + 0.3148 * \text{최고기온} - 0.1144 * 99\text{강수} + 0.2207 * \text{합계일조시간} - 0.0644 * \text{초미세먼지}$

$- 0.0607 * \text{일강수량} + 0.6322 * \text{미세먼지_기준_2} + 1.2346 * \text{미세먼지_기준_3} - 94.5226 * \text{미세먼지_기준_4} - 1.0285 * \text{휴일여부_1}$

Conclusion

Contents

Logistic Regression
/ Model Assessment

결론

따릉이 이용객 수를 평균값을 기준으로 하여 많다 / 적다로 구분한 다음, 날씨 데이터를 이용한 logistic regression 모델로 따릉이 이용객 수를 예측함.

전체 데이터 중 5개의 숫자형 변수(최고기온, 9-9강수, 합계일조시간, 초미세먼지, 일강수량)과 2개의 범주형 변수(미세먼지_기준, 휴일여부)를 전진 선택법을 활용하여 선택한 다음 VIF로 다중공선성을 확인함. 이 과정에서 이상치의 존재 여부를 확인했고 outlier가 $Q3 + 1.5IQR$ 안으로 들어오도록 보정 진행.

모델 성능 평가 결과,

정확도 = 88.14%, 특이도 = 78.57%, 정밀도 = 83.33%, 재현율 = 96.77%, F1 score = 89.55%로 종합적으로 높은 성능을 보임.

한계

Logistic Regression을 위한 6가지 가정 중 일부 조건은 만족하였으나 충족되지 못한 조건들이 존재함.

- 날씨 데이터와 따릉이 데이터를 사용했는데, 이는 날짜에 따른 데이터이기 때문에 시계열 데이터임. 따라서 관측치들 간에 독립성이 성립하지 않음.
- 독립 변수와 log odds 간의 linearity를 확인하였을 때 최고기온처럼 선형성을 보이는 데이터도 있었으나 대부분 선형성을 보이지 않았음.
- 데이터들 간의 등분산성이 성립되지 않았음.

개선점

Data 관련 개선점

관측치의 독립성 : 시계열 데이터가 아닌 반복적 측정 X , 시간 및 공간에 따른 상관관계 X 데이터 활용을 통해 관측치들 간 독립성 확보 가능

변수와 log odds의 선형성 : polynomial terms 등 고차 다항식 항식을 통합하거나 연속형을 범주형으로 변환 통해 독립 변수의 선형 패턴을 확보함으로써 로짓 선형성 가정을 만족할 수 있음

Data Source

Contents

Logistic Regression
/ Model Assessment

[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>

(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>