

제주도 호텔 가격 회귀 분석

4조 이진서, 김형기





01 프로젝트 목적 및 동기

데이터 수집

02

crawling, API 이용

데이터 전처리 (EDA)

03

성급, 리뷰, 평점 가격 분석

회귀 분석

04

분석 및 모델 평가

결론 및 한계점

05

예측 값 확인 및 한계점



01 프로젝트 목적 및 동기



02 데이터 수집

1. 인터파크 Crawling
2. Booking.com API

인터파크 Crawling

```
1 # 함수로 만들기
2 from selenium import webdriver
3 from time import sleep
4 from datetime import datetime
5 import time
6
7 datas = []
8
9 def get_hotels(cin):
10     url = 'https://travel.interpark.com/checkinnow/search?\
11 q=%EC%A0%9C%EC%A3%BC,%20%EC%A0%84%EC%B2%B4&disp_q=%EC%A0%9C%EC%A3%BC,\\
12 %20%EC%A0%84%EC%B2%B4&coord=&regionidisp=42090&startdate={}&enddate={}&\n\
13 roomOptions=0%5E2%5E0%5Ecategory1=%EC%88%99%EB%B0%95&type=&value=&index=1\
14 &npprmaxsort=min'.format(cin, cin+1)
15     driver = webdriver.Chrome()
16     driver.get(url)
17
18     # 숙소 유형 호텔 선택
19     driver.find_element_by_xpath('//*[@@id="resultContainer"]/div[1]/div/div/dl[2]/dd[1]/label[2]').click()
20     time.sleep(3)
21
22     # 1, 2, 3, 4, 5 등급 선택
23     driver.find_element_by_xpath('//*[@@id="resultContainer"]/div[1]/div/div/dl[3]/dd/label[1]').click()
24     driver.find_element_by_xpath('//*[@@id="resultContainer"]/div[1]/div/div/dl[3]/dd/label[2]').click()
25     driver.find_element_by_xpath('//*[@@id="resultContainer"]/div[1]/div/div/dl[3]/dd/label[3]').click()
26     driver.find_element_by_xpath('//*[@@id="resultContainer"]/div[1]/div/div/dl[3]/dd/label[4]').click()
27     driver.find_element_by_xpath('//*[@@id="resultContainer"]/div[1]/div/div/dl[3]/dd/label[5]').click()
28     time.sleep(3)
29
30     # 총 페이지 갯수
31     total_page = (int(driver.find_element_by_css_selector("#sortMenuCount").text)//30)+1
32
33     for page in range(1, total_page+1):
34         try:
35             print(page)
36             driver.execute_script("window.scrollTo(0, document.body.scrollHeight);") # 스크롤 내리고 로딩 기다리기
37             time.sleep(3)
38
39             hlist = driver.find_elements_by_xpath("//div[contains(@class, 'col data')]")
40
41             for i in range(len(hlist)):
42                 name = hlist[i].find_element_by_css_selector("div:nth-child(1) > div.titArea > a > strong").text
43                 price = hlist[i].find_element_by_css_selector("div:nth-child(1) > div.priceArea > strong").text
44                 star = hlist[i].find_element_by_css_selector('div:nth-child(2) > div > span.inline > span:nth-c\
45 r_count= hlist[i].find_element_by_css_selector('div:nth-child(3) > div.ratingArea > span.review\
46 addr = hlist[i].find_element_by_css_selector('div:nth-child(2) > div > span:nth-child(1)').text
47                 try:
48                     r_score = hlist[i].find_element_by_css_selector('div:nth-child(3) > div.ratingArea > span.s\
49                 except:
50                     r_score = 0
51
52                 datas.append({"cin":cin, "name":name, "price":price, "star":star,
53                               "r_count":r_count, "r_score":r_score, "addr":addr})
54
55             page_selector = ".paging > span > a:nth-child({})".format(str(page+1))
56             driver.find_element_by_css_selector(page_selector).click() # 다음페이지 클릭
57             time.sleep(3)
58         except:
59             break
60
61
62         print("crawling : {}".format(cin))
63         driver.quit()
```

- Check_in_date
- Name
- Price
- Star
- Rating
- Reviews

booking.com Crawling

Extract data from Booking.com where there is no official API. With one click, you can scrape all information about accommodation listings, such as price, rating, reviews, stars, etc. It's easy to use and input is similar to the website.

Try for free

Editor JSON

Booking search query: Jeju Island

Destination type: City

Start URLs: + Add URL, Link remote text file, Upload text file

Order results by: Stars [5->1]

Minimum rating:

Maximum pagination pages:

Check-in date (yyyy-mm-dd): 2020-09-29

Check-out date (yyyy-mm-dd): 2020-09-30

Number of rooms: 1

Number of adults: 2

Number of children: 0

Preferred currency: KRW - Korean won

Language: 한국어

Price range: none

- *Check_in_date*
- *Name*
- *Price*
- *Features*
- *Rating*
- *Reviews*



03 데이터 전처리

1. 인터파크 EDA
2. Booking.com EDA

인터파크 EDA

```
1 # 경도 위도 naver API
2 import numpy as np
3 import pandas as pd
4 from urllib.request import urlopen
5 from urllib import parse
6 from urllib.request import Request
7 from urllib.error import HTTPError
8 from bs4 import BeautifulSoup
9 import json
10
11 # naver api
12 client_id = '0d9uya4hfz';      # 할당받은 ID
13 client_pw = 'zB1NWFxWwySehisDj6ey0lu2EJ33HgtPagEwCaU';      # 할당받은 Secret
14
15 api_url = 'https://naveropenapi.apigw.ntruss.com/map-geocode/v2/geocode?query='
16
17 # 주소 목록 파일
18 data = hotels_names
19 data
20
21 # 네이버 지도 API 이용해서 위경도 찾기
22 geo_coordi = []
23 for addr in hotels_names['addr']:
24     addr_urllenc = parse.quote(addr)
25     url = api_url + addr_urllenc
26     request = Request(url)
27     request.add_header('X-NCP-APIGW-API-KEY-ID', client_id)
28     request.add_header('X-NCP-APIGW-API-KEY', client_pw)
29     try:
30         response = urlopen(request)
31     except HTTPError as e:
32         print('HTTP Error!')
33         latitude = None
34         longitude = None
35     else:
36         rescode = response.getcode()
37         if rescode == 200:
38             response_body = response.read().decode('utf-8')
39             response_body = json.loads(response_body) # json
40             if response_body['addresses'] == []:
41                 print("result' not exist!")
42                 latitude = None
43                 longitude = None
44             else:
45                 latitude = response_body['addresses'][0]['y']
46                 longitude = response_body['addresses'][0]['x']
47                 print("Success!")
48         else:
49             print('Response error code : %d' % rescode)
50             latitude = None
51             longitude = None
52
53     geo_coordi.append([latitude, longitude])
54
55
56 np_geo_coordi = np.array(geo_coordi)
57 pd_geo_coordi = pd.DataFrame({'name': hotels_names['name'].values,
58                               'addr': hotels_names['addr'].values,
59                               'latitude': np_geo_coordi[:, 0],
60                               'longitude': np_geo_coordi[:, 1]})
```

```
In [119]: 1 # 제주공항과 거리는? Haversine
2 !pip install haversine
Collecting haversine
  Downloading haversine-2.2.0-py2.py3-none-any.whl (4.9 kB)
Installing collected packages: haversine
Successfully installed haversine-2.2.0

In [4]: 1 # 제주공항과 거리는?
2 from haversine import haversine
In [5]: 1 len(lnglat_df)
Out[5]: 75

In [6]: 1 distance_list = []
2 data1 = (33.511005, 126.491318) # 제주공항
3
4 for i in range(len(lnglat_df)):
5     data2 = (lnglat_df["latitude"][i], lnglat_df["longitude"][i]) # 각 호텔들의 위경도
6     distance = haversine(data1, data2, unit ='km') # 제주공항과 호텔의 거리
7     distance_list.append(distance)
8
9
10 lnglat_df['distance'] = distance_list
11 lnglat_df.drop(columns=['name'], axis=1, inplace=True)
12 lnglat_df

Out[6]:
```

	addr	latitude	longitude	distance
0	제주 서귀포시 안덕면 서귀리 산 24번지	33.307187	126.317436	27.822325
1	제주특별자치도 제주시 조천읍 함덕리 1269-9 유립유블레스호텔	33.542411	126.666212	16.583734
2	제주 제주시 건입동 1443	33.517411	126.526897	3.374514
3	제주 서귀포시 색달동 2812-4	33.248275	126.410407	30.164860
4	제주 제주시 이도이동 1025-8	33.493516	126.534292	4.433829
...
70	제주 서귀포시 표선면 표선리 40-69	33.323364	126.844747	38.876091
71	제주 제주시 외도동 401-8	33.492431	126.428735	6.159437
72	제주 서귀포시 회수동 30	33.286153	126.444396	25.379001
73	제주 서귀포시 색달동 2138	33.259434	126.406306	29.065655
74	제주 서귀포시 색달동 3039-3	33.247167	126.407950	30.341594

75 rows × 4 columns

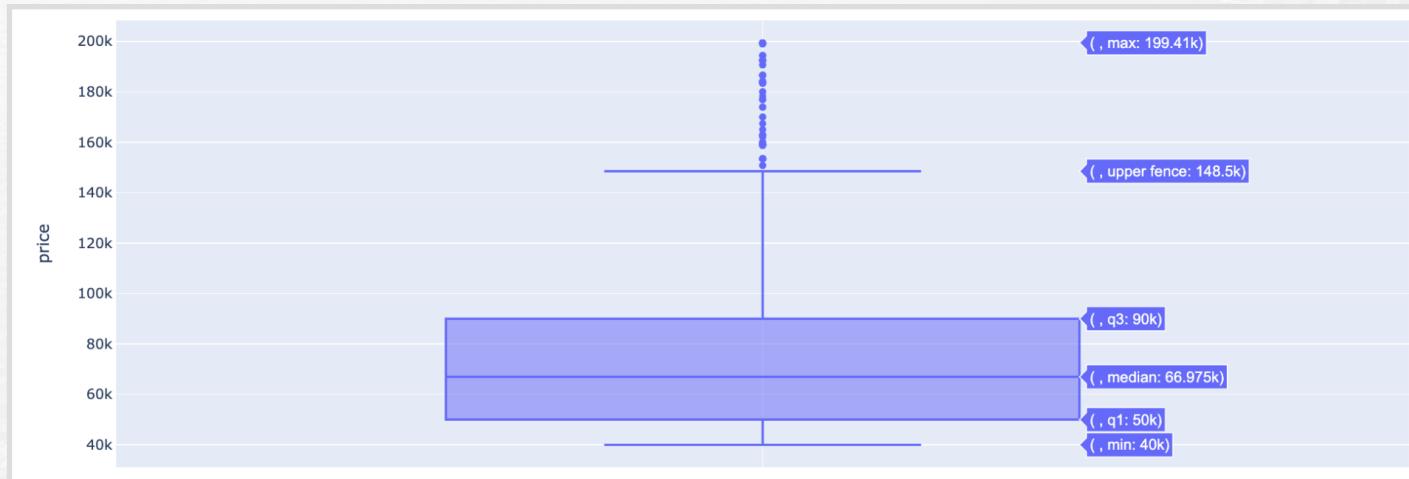
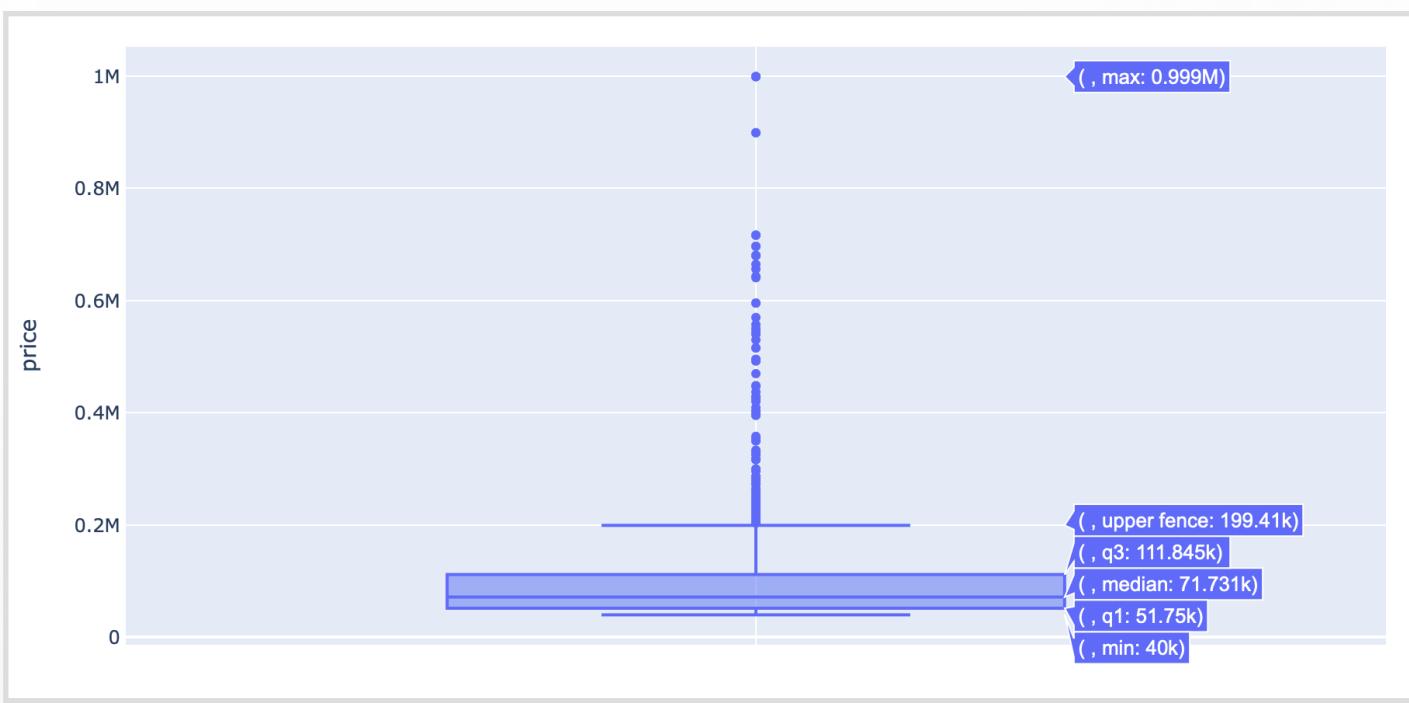
인터파크 EDA

In [8]:	1	hotels_df
Out [8]:		
		cin name price star r_count r_score addr latitude longitude distance
0	20200901	메리어트관 제주신화월드 호텔앤리조트 353760 5 224 9.0 제주 서귀포시 안덕면 서광리 산 24번지 33.307187 126.317436 27.823235
1	20200902	메리어트관 제주신화월드 호텔앤리조트 353760 5 224 9.0 제주 서귀포시 안덕면 서광리 산 24번지 33.307187 126.317436 27.823235
2	20200903	메리어트관 제주신화월드 호텔앤리조트 353760 5 224 9.0 제주 서귀포시 안덕면 서광리 산 24번지 33.307187 126.317436 27.823235
3	20200904	메리어트관 제주신화월드 호텔앤리조트 353760 5 224 9.0 제주 서귀포시 안덕면 서광리 산 24번지 33.307187 126.317436 27.823235
4	20200905	메리어트관 제주신화월드 호텔앤리조트 316910 5 224 9.0 제주 서귀포시 안덕면 서광리 산 24번지 33.307187 126.317436 27.823235
...
1917	20200925	제주신라호텔 595500 5 4863 9.8 제주 서귀포시 색달동 3039-3 33.247167 126.407950 30.341594
1918	20200926	제주신라호텔 544500 5 4863 9.8 제주 서귀포시 색달동 3039-3 33.247167 126.407950 30.341594
1919	20200927	제주신라호텔 656000 5 4863 9.8 제주 서귀포시 색달동 3039-3 33.247167 126.407950 30.341594
1920	20200928	제주신라호텔 680200 5 4863 9.8 제주 서귀포시 색달동 3039-3 33.247167 126.407950 30.341594
1921	20200929	제주신라호텔 680200 5 4863 9.8 제주 서귀포시 색달동 3039-3 33.247167 126.407950 30.341594

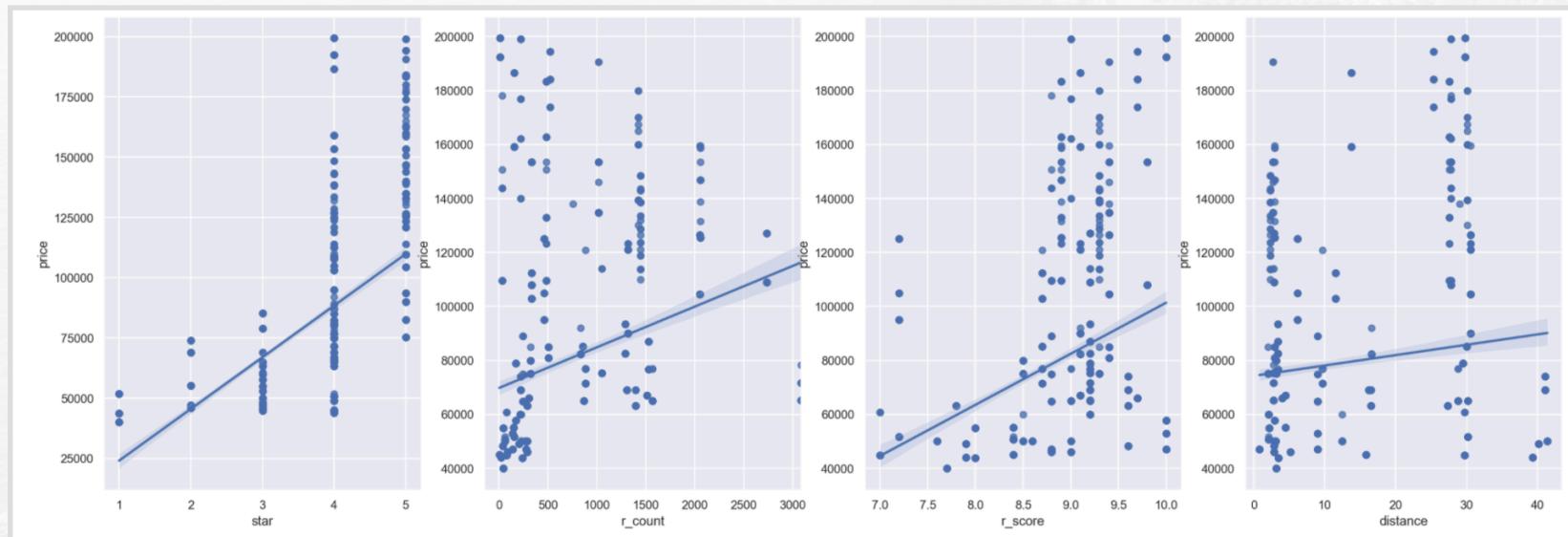
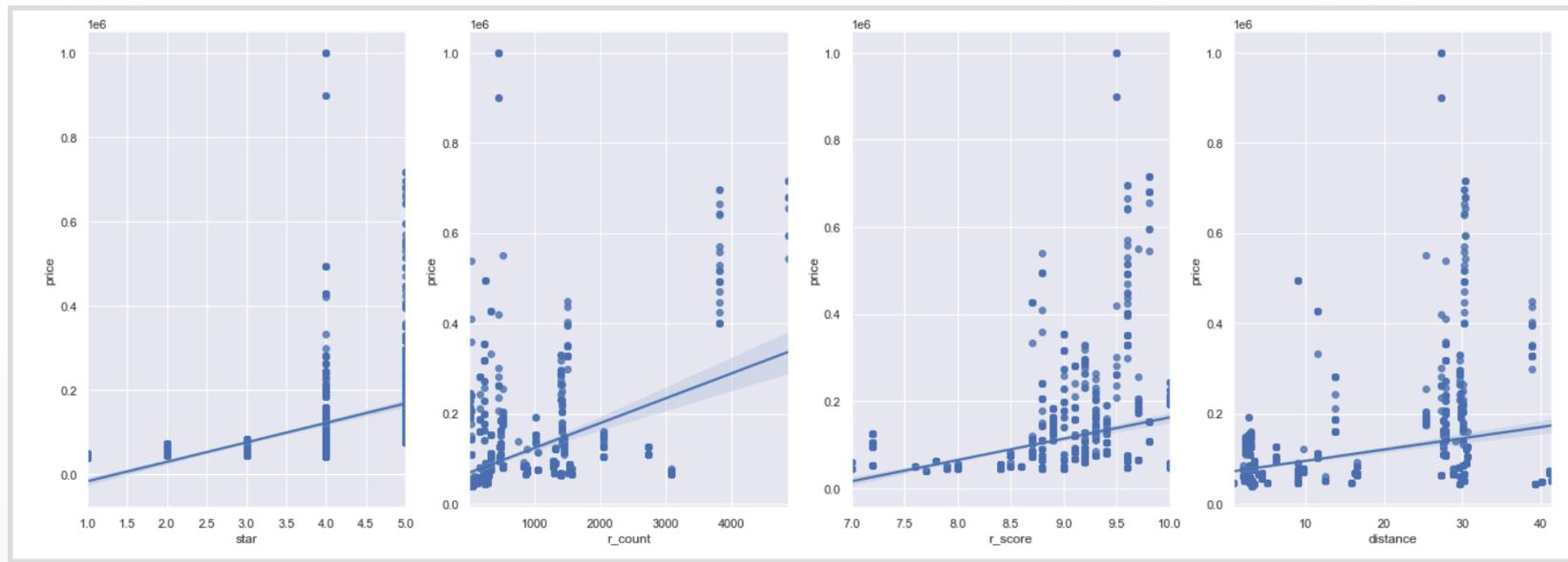
1922 rows × 10 columns

- *Check_in_date(cin)*
- *Name*
- *Price*
- *Star*
- *Rating (r_score)*
- *Reviews (r_count)*
- *Distance*

인터파크 EDA



인터파크 EDA



Booking.com EDA

```
In [56]: 1 df = pd.read_csv('/Users/stevenkim/Downloads/booking.com 원본/원본/dataset_CHI_2020_09_01.csv')
2 df
```

270	대한민국	제주, 281, Hamwa-ro, Jocheon-eup	63334	제주도	281, Hamwa-ro, Jocheon-eup	NaN	NaN	NaN	제주 베스트셀러! 제주에 위치한 힐링 예쁜 집게스트하우스는 무료 WiFi, 에어컨 ...	https://cf.bstatic.com/ir
271	대한민국	제주, 281, Hamwa-ro, Jocheon-eup	63334	제주도	281, Hamwa-ro, Jocheon-eup	NaN	NaN	NaN	제주 베스트셀러! 제주에 위치한 힐링 예쁜 집게스트하우스는 무료 WiFi, 에어컨 ...	https://cf.bstatic.com/ir
272	대한민국	제주, 남성로 131	63172	제주도	남성로 131	유럽식	15:00	20:00	참피온 호텔은 제주 국제공항에서 동쪽으로 2.5km 떨어져 있습니다. 인기 관광지인...	https://cf.bstatic.com/ir

273 rows > 2967 columns

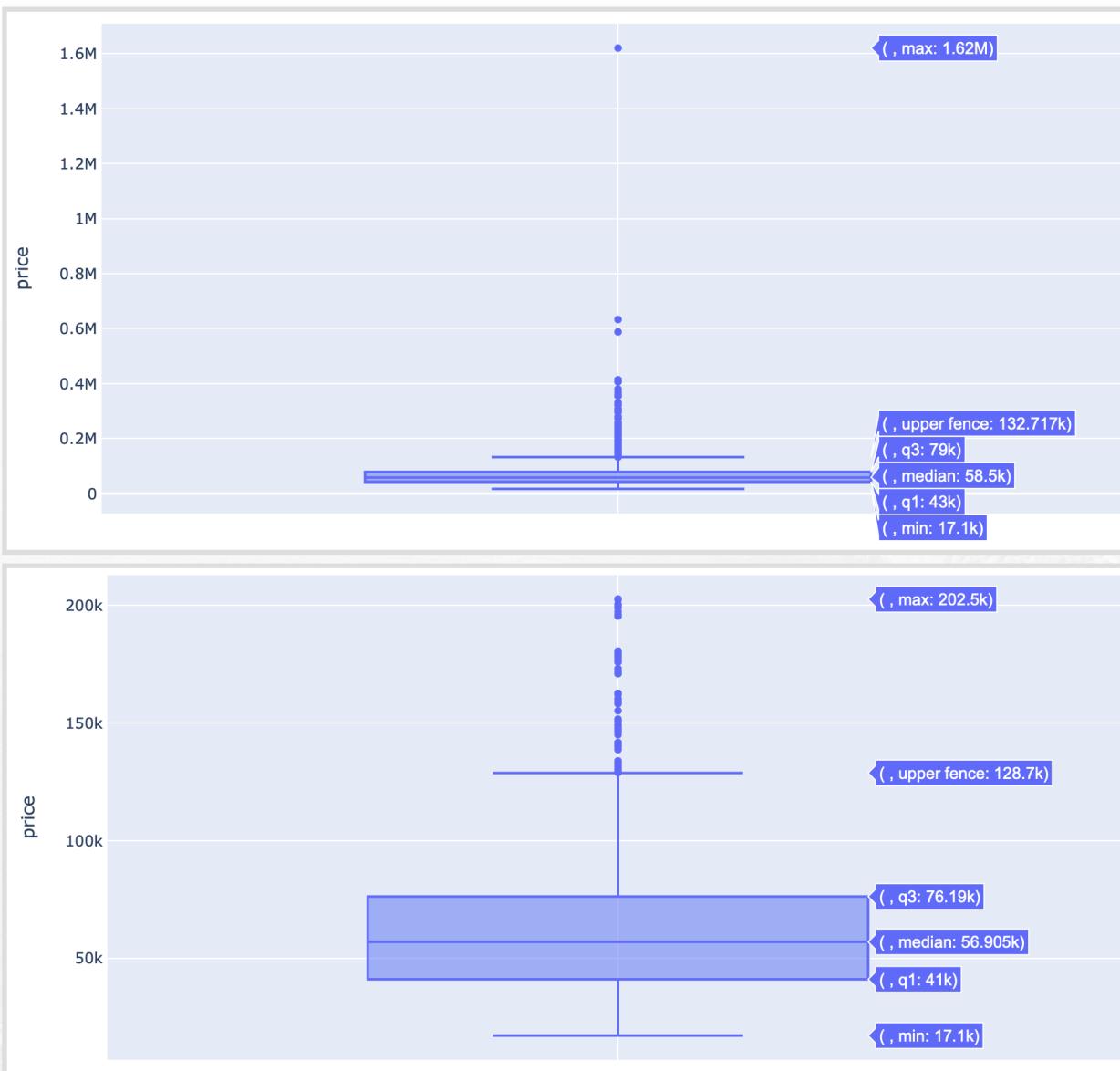
```
In [3]: 1 hotels_df = pd.concat([df02, df03, df04, df05, df06, df07, df08, df09, df10, df11, df12, df13, df14, df15, df16,
2 hotels_df = hotels_df[['name', 'price', 'Area', 'rating', 'reviews', 'distance', 'day', 'weekend']]
3 hotels_df
```

Out [3]:

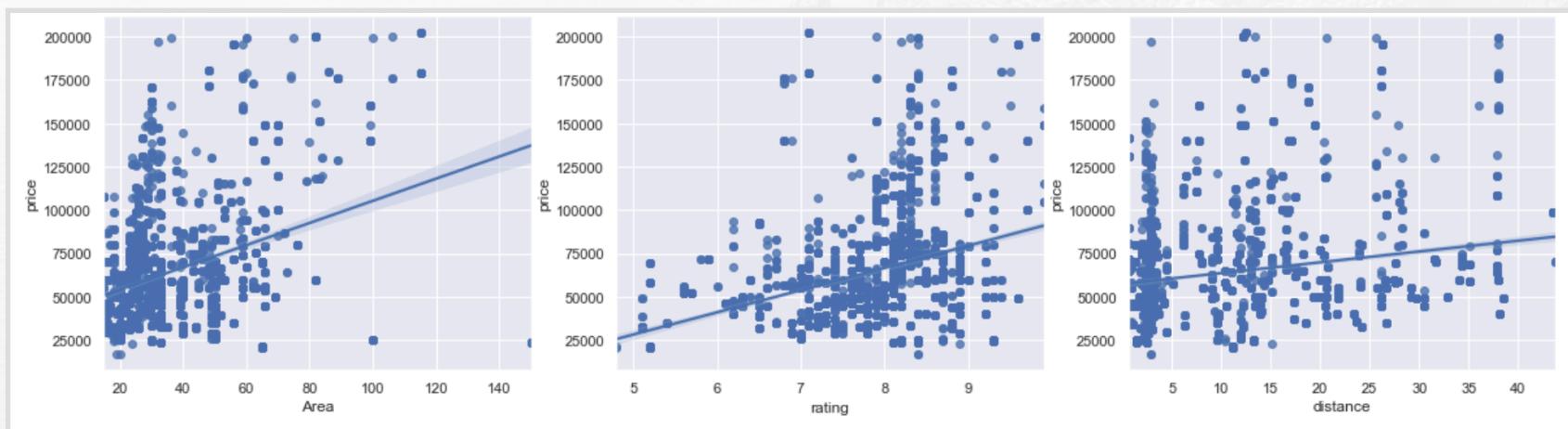
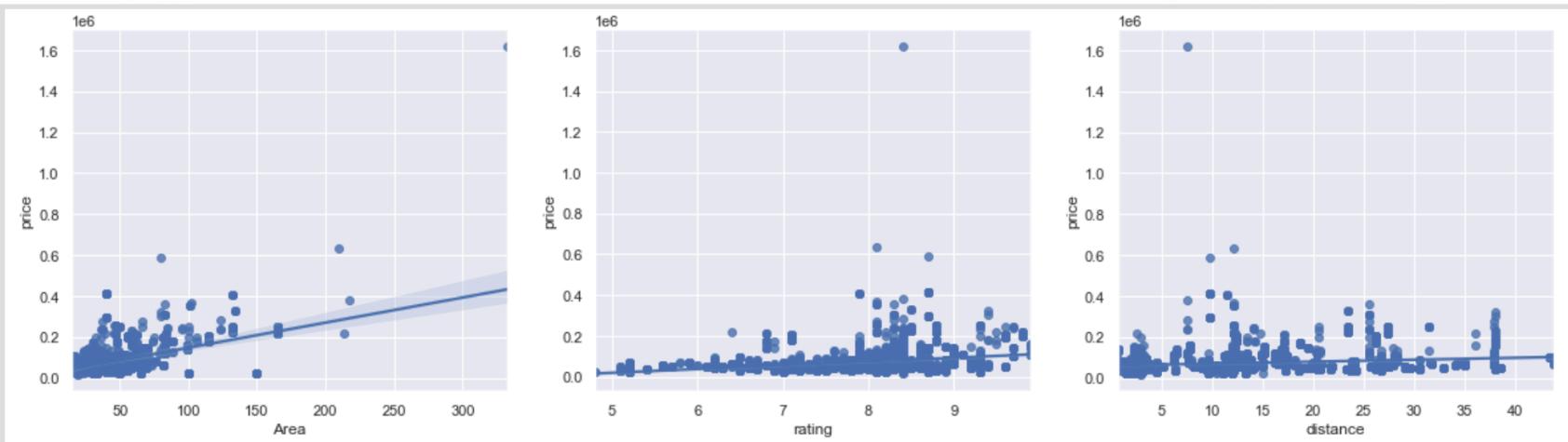
	name	price	Area	rating	reviews	distance	day	weekend
0	제주오리엔탈호텔 & 카지노	67725	24	7.7	587	2.738362	20200902	0
1	다인 오세아노 호텔	355555	101	8.1	541	12.129263	20200902	0
2	블랙 샌즈 호텔 제주	49600	25	7.7	120	8.986851	20200902	0
3	제주 베스트 힐	79000	59	9.0	20	16.686766	20200902	0
4	제주 하바나 리조트	59000	36	7.7	230	35.060464	20200902	0
...
4899	제주 삼다 호스텔	46540	45	7.6	9	0.830589	20200928	0
4900	예스준 게스트하우스	25000	20	8.7	12	1.922235	20200928	0
4901	미르 게스트하우스	58500	23	8.5	305	1.809006	20200928	0
4902	예하 게스트하우스	60000	23	8.9	545	2.748919	20200928	0
4903	호텔 메르블루	130784	27	8.7	8	0.657127	20200928	0

4904 rows > 8 columns

Booking.com EDA



Booking.com EDA



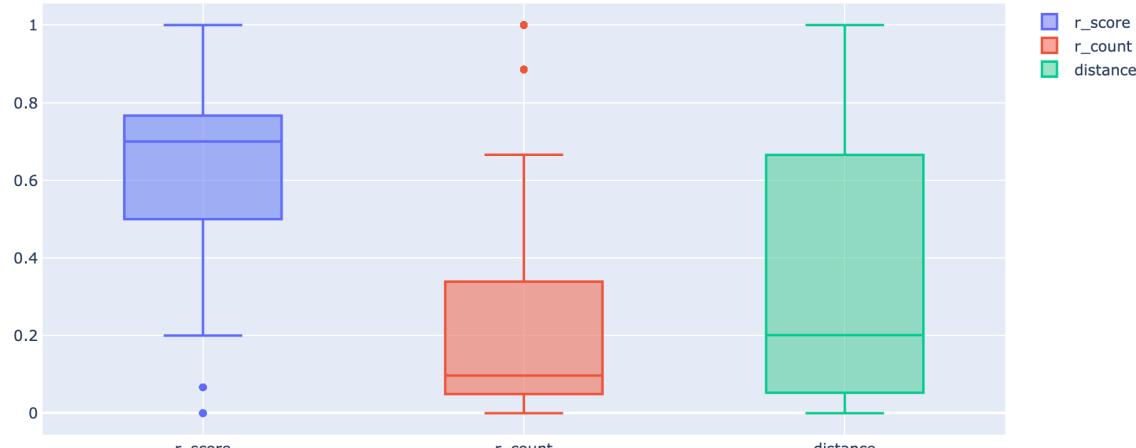


04 회귀 분석

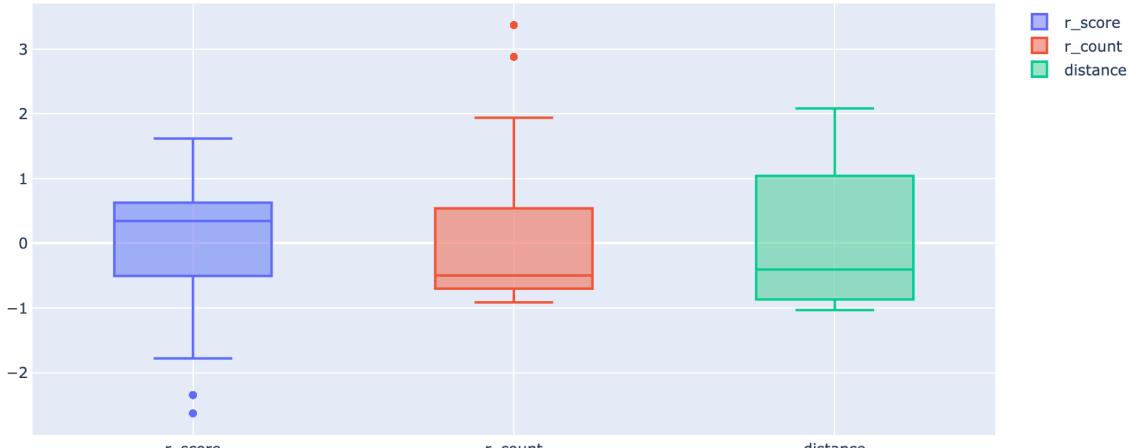
- 1. 인터파크 회귀 분석
- 2. Booking.com 회귀 분석

인터파크 회귀 분석

MinMaxScaler



StandardScaler



인터파크 회귀 분석

```
In [69]: 1 X = hotels_df_copy[["cin", "star", "r_count", "r_score", "distance"]]
2 y = hotels_df_copy["price"]
```

```
In [70]: 1 from sklearn.model_selection import train_test_split
2
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=7)
```

```
In [71]: 1 import numpy as np
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error, r2_score
4
5 lr = LinearRegression()
6 lr.fit(X_train, y_train)
7
8 pred_tr = lr.predict(X_train)
9 pred_test = lr.predict(X_test)
10
11 rmse_tr = (np.sqrt(mean_squared_error(y_train, pred_tr)))
12 rmse_test = (np.sqrt(mean_squared_error(y_test, pred_test)))
13
14 print("RMSE of Train Data : ", rmse_tr)
15 print("RMSE of Test Data : ", rmse_test)
16
17 print("r2 score of Train Data : ", r2_score(y_train, pred_tr))
18 print("r2 score of Test Data : ", r2_score(y_test, pred_test))
```

```
RMSE of Train Data :  28579.599102446613
RMSE of Test Data :  28703.443898098023
r2 score of Train Data :  0.4069110283918662
r2 score of Test Data :  0.39393124632773013
```

인터파크 회귀 분석

```
In [75]: 1 models = []
2 models.append('LinearReg', LinearRegression())
3 models.append('DecisionTree', DecisionTreeRegressor())
4 models.append('RandomForest', RandomForestRegressor())
5 models.append('GradientB', GradientBoostingRegressor())
6 models.append('XGradianB', XGBRegressor())
7 models.append('LightGBM', LGBMRegressor())
```

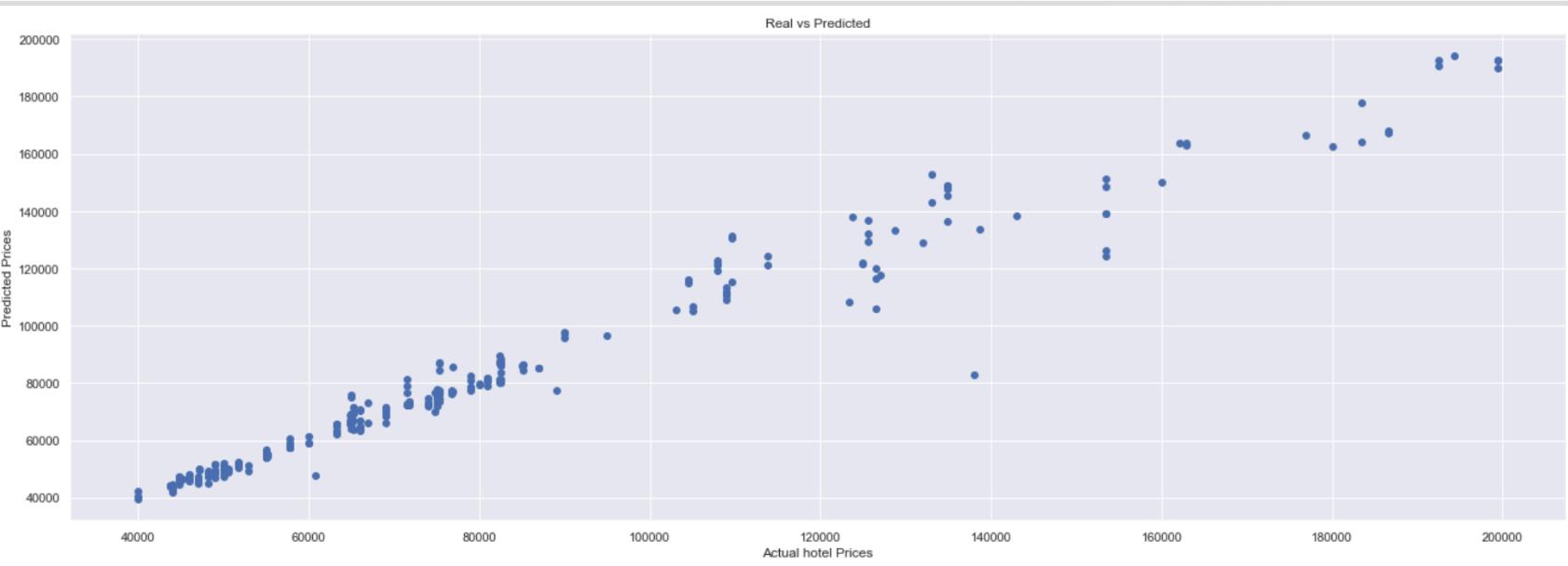
```
In [76]: 1 names = []
2 train_rmse = []
3 test_rmse = []
4 train_r2 = []
5 test_r2 = []
6
7 for name, model in models:
8     model.fit(X_train, y_train)
9     pred_test = model.predict(X_test)
10    pred_train = model.predict(X_train)
11
12    names.append(name)
13    train_rmse.append((np.sqrt(mean_squared_error(y_train, pred_train))))
14    test_rmse.append((np.sqrt(mean_squared_error(y_test, pred_test))))
15    test_r2.append(r2_score(y_test, pred_test))
16    train_r2.append(r2_score(y_train, pred_train))
```

```
In [77]: 1 result = pd.DataFrame({'model name':names,
2                           'train rmse':train_rmse,
3                           'test rmse':test_rmse,
4                           'train score':train_r2,
5                           'test score':test_r2})
6 result
```

Out[77]:

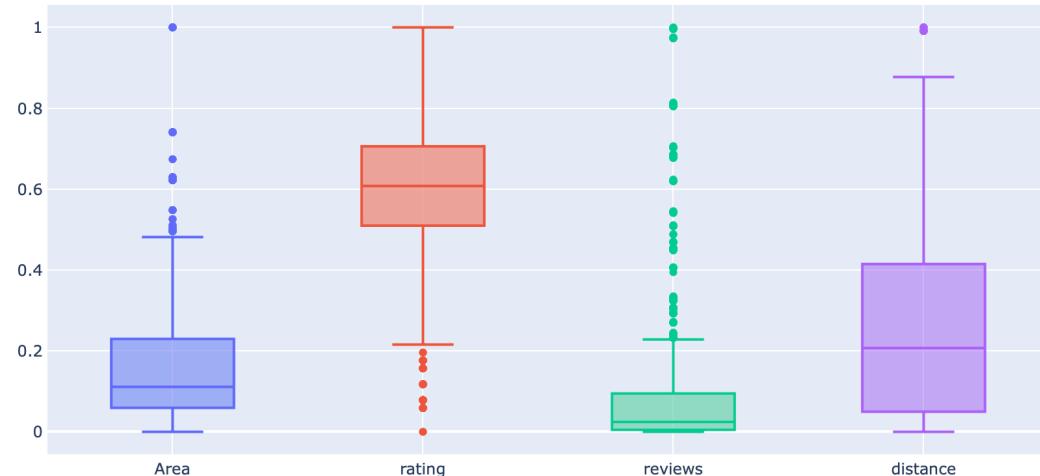
	model name	train rmse	test rmse	train score	test score
0	LinearReg	28579.599102	28703.443898	0.406911	0.393931
1	DecisionTree	4363.608472	8081.181007	0.986174	0.951960
2	RandomForest	4754.385280	7377.912814	0.983587	0.959958
3	GradientB	8100.342462	8627.853661	0.952355	0.945241
4	XGradianB	4382.047461	7868.607918	0.986057	0.954454
5	LightGBM	5043.809811	6861.865932	0.981528	0.965363

인터파크 회귀 분석

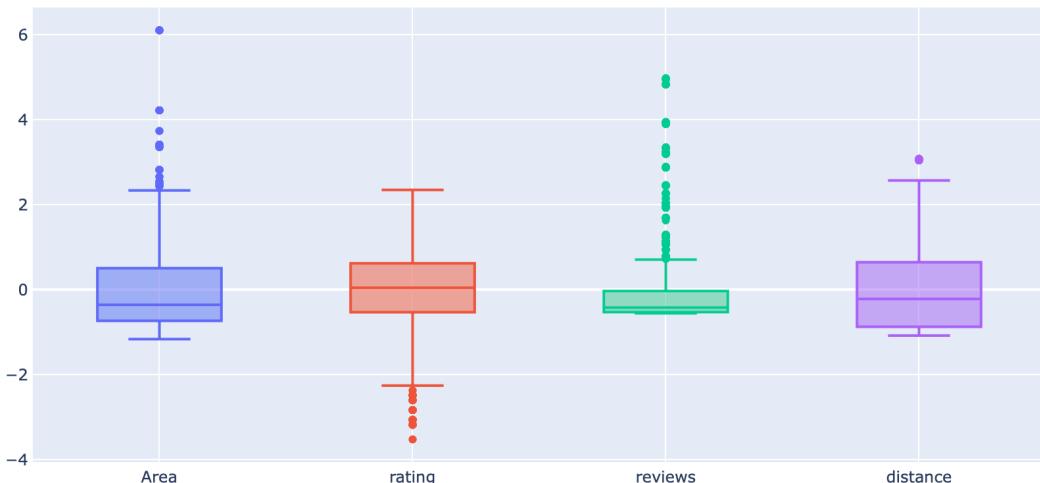


Booking.com 회귀 분석

MinMaxScaler



StandardScaler



Booking.com 회귀 분석

LinearRegression

In [18]:

```
1 from sklearn.model_selection import train_test_split
2
3 X = hotels_df_copy[['Area', 'rating', 'reviews', 'distance', 'weekend']]
4 y = hotels_df_copy['price']
5
6 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=13)
```

In [26]:

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import r2_score, mean_squared_error
3
4 lr = LinearRegression()
5 lr.fit(X_train, y_train)
6
7 pred_tr = lr.predict(X_train)
8 pred_test = lr.predict(X_test)
9
10 rmse_tr = (np.sqrt(mean_squared_error(y_train, pred_tr)))
11 rmse_test = (np.sqrt(mean_squared_error(y_test, pred_test)))
12
13 print("RMSE of Train Data : ", rmse_tr)
14 print("RMSE of Test Data : ", rmse_test)
15
16 print('Train Acc : ', r2_score(y_train, pred_tr))
17 print('Test Acc : ', r2_score(y_test, pred_test))
```

```
RMSE of Train Data :  28706.68841444606
RMSE of Test Data :  27549.056886837352
Train Acc :  0.2713355744005723
Test Acc :  0.2962487550499314
```

Booking.com 회귀 분석

```
In [70]: 1 models = []
2 models.append(('LinearReg', LinearRegression()))
3 models.append(('DecisionTree', DecisionTreeRegressor()))
4 models.append(('RandomForest', RandomForestRegressor()))
5 models.append(('GradientB', GradientBoostingRegressor()))
6 models.append(('XGradianB', XGBRegressor()))
7 models.append(('LightGBM', LGBMRegressor()))
```

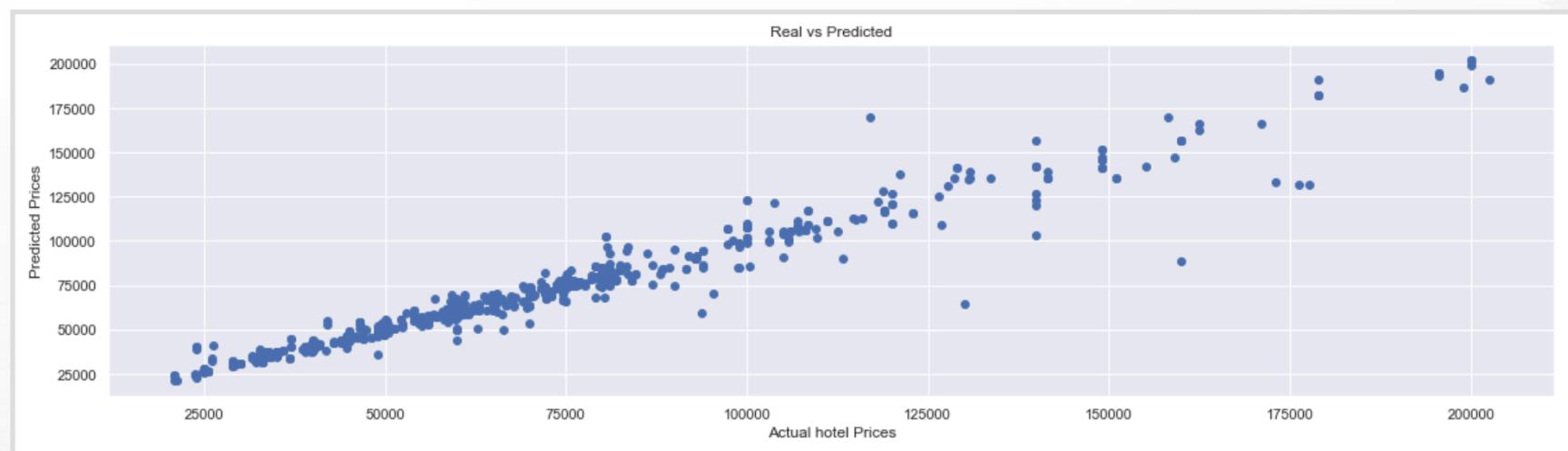
```
In [71]: 1 names = []
2 train_rmse = []
3 test_rmse = []
4 train_r2 = []
5 test_r2 = []
6
7 for name, model in models:
8     model.fit(X_train, y_train)
9     pred_test = model.predict(X_test)
10    pred_train = model.predict(X_train)
11
12    names.append(name)
13    train_rmse.append((np.sqrt(mean_squared_error(y_train, pred_train))))
14    test_rmse.append((np.sqrt(mean_squared_error(y_test, pred_test))))
15    test_r2.append(r2_score(y_test, pred_test))
16    train_r2.append(r2_score(y_train, pred_train))
```

```
In [72]: 1 result = pd.DataFrame({'model name':names,
2                           'train rmse':train_rmse,
3                           'test rmse':test_rmse,
4                           'train score':train_r2,
5                           'test score':test_r2})
6 result
```

Out[72]:

	model name	train rmse	test rmse	train score	test score
0	LinearReg	28706.688414	27549.056887	0.271336	0.296249
1	DecisionTree	2976.038312	6062.792312	0.992169	0.965916
2	RandomForest	3378.216818	5348.107250	0.989909	0.973478
3	GradientB	13484.785516	14077.261969	0.839214	0.816244
4	XGradianB	3074.690890	5495.817758	0.991641	0.971993
5	LightGBM	4882.735552	6539.204679	0.978919	0.960349

Booking.com 회귀 분석





05 결론 및 한계점

1. 예측 값 확인
2. 회고 (한계점)

예측 값 확인

인터파크

```
In [59]: 1 jinseo = np.array([[20200925, 4, 1000, 8.5, 5]])
2 print("Jinseo : ", best_model.best_estimator_.predict(jinseo))
```

```
Jinseo : [71990.44158191]
```

```
In [60]: 1 hyunggi = np.array([[20200918, 4, 750, 8, 10]])
2 print("Hyunggi : ", best_model.best_estimator_.predict(hyunggi))
```

```
Hyunggi : [94817.71601341]
```

booking.com

```
In [83]: 1 # Area, Reviews, rating, distance, weekend (0:weekday)
2 jinseo = np.array([[38, 243, 7.8, 12, 0]])
3 print("Jinseo : ", best_model_ss.best_estimator_.predict(jinseo))
```

```
Jinseo : [140252.08570413]
```

```
In [84]: 1 # Area, Reviews, rating, distance, weekend (1:weekend)
2 hyunggi = np.array([[38, 243, 7.8, 12, 1]])
3 print("Hyunggi : ", best_model_ss.best_estimator_.predict(hyunggi))
```

```
Hyunggi : [141702.38401301]
```

한계점

가격 차이

조회기간 : 2019.07 ~ 2019.07, 여객 : 전체, 화물 : 전체, 공항 : 제주(CJU), 운항 : 전체, 노선 : 전체, 여객화물 : 전체				
공항명	공급(석)	운항(편)	여객(명)	
제주(CJU)	1,619,236	8,338	1,404,339	

조회기간 : 2020.07 ~ 2020.07, 여객 : 전체, 화물 : 전체, 공항 : 제주(CJU), 운항 : 전체, 노선 : 전체, 여객화물 : 전체				
공항명	공급(석)	운항(편)	여객(명)	
제주(CJU)	1,236,092	6,405	980,222	

(source : 항공정보포털시스템 (www.airportal.go.kr))

대략 30 % 감소



수요 감소가 주중, 주말 가격 형성에 영향 미쳤는지?

크롤링

데이터 수집에 제약

→ 데이터 컬럼들이 많았으면, 다양한 요인들을 바탕으로 가격 분석을 진행할 수 있지 않았을까?

THANK YOU

