

# Análisis Predictivo

## TP2

Profesores:

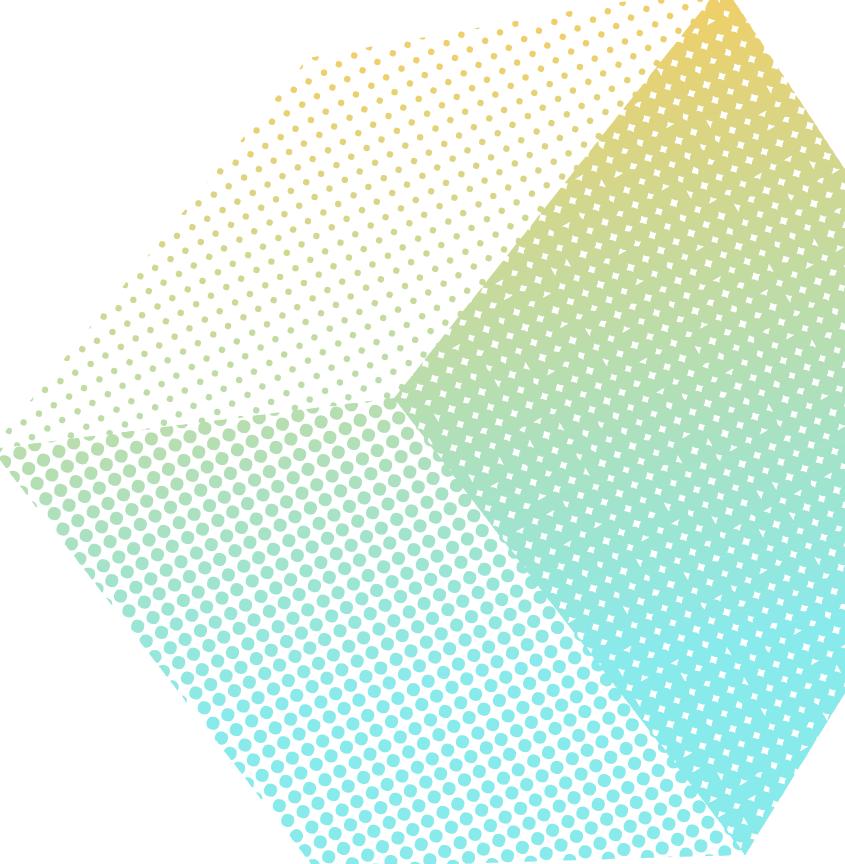
- Ezequiel Martín Eliano Sombory
- Leonardo Andrés Caravaggio
- Francisco Valentin

Nicolás Rodrigues da Cruz



ITBA

# Score obtenido en Kaggle



|   |   |                           |   |         |    |    |
|---|---|---------------------------|---|---------|----|----|
| 8 | — | NICOLAS RODRIGUES DA CRUZ |  | 0.68842 | 22 | 5d |
|---|---|---------------------------|---|---------|----|----|



# Índice de contenidos

**01**

VISTAZO GENERAL

**02**

TRATAMIENTO DE VALORES NULOS Y OUTLIERS

**03**

VARIABLE OBJETIVO

**04**

ENCODING DE VARIABLES Y CREACION DE NUEVAS

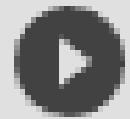
**05**

MODELOS Y PREDICCIONES

# 01. Vistazo general

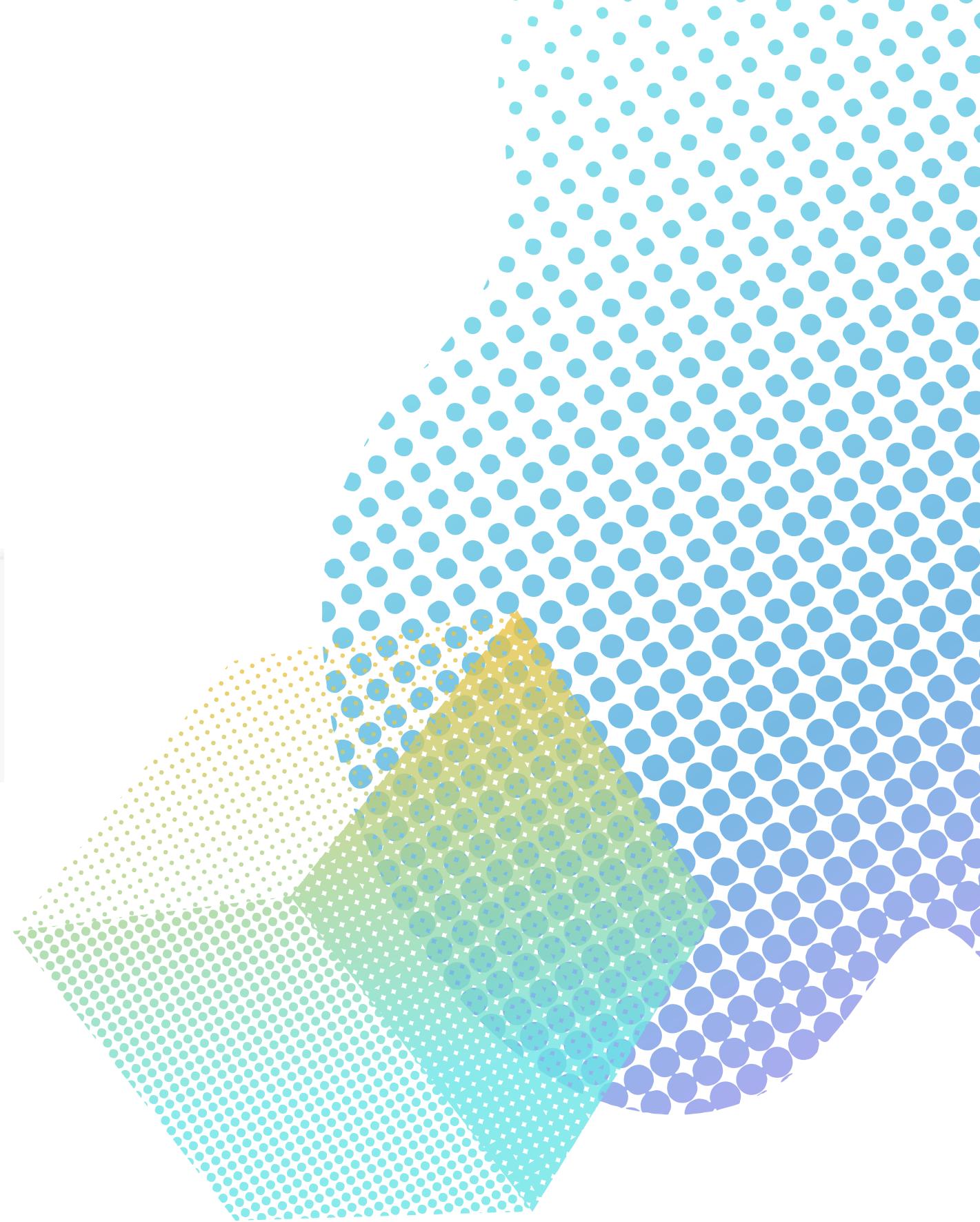
| Características | Cantidad |
|-----------------|----------|
| Filas           | 4928     |
| Columnas        | 68       |
| Duplicados      | 0        |

# 01. Vistazo general

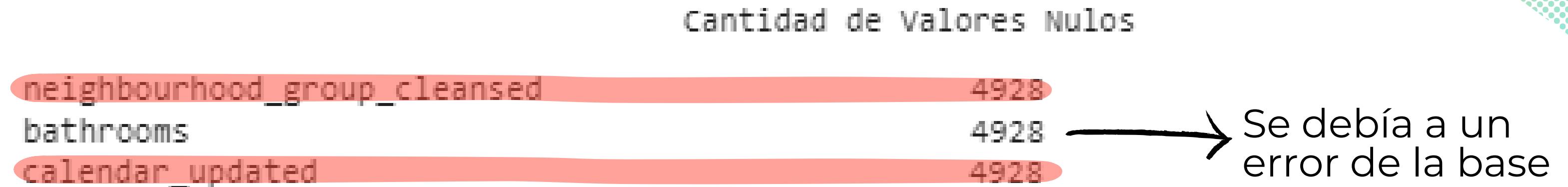


```
data_types = train.dtypes.value_counts()  
print(data_types)
```

|         |    |
|---------|----|
| object  | 29 |
| int64   | 22 |
| float64 | 17 |



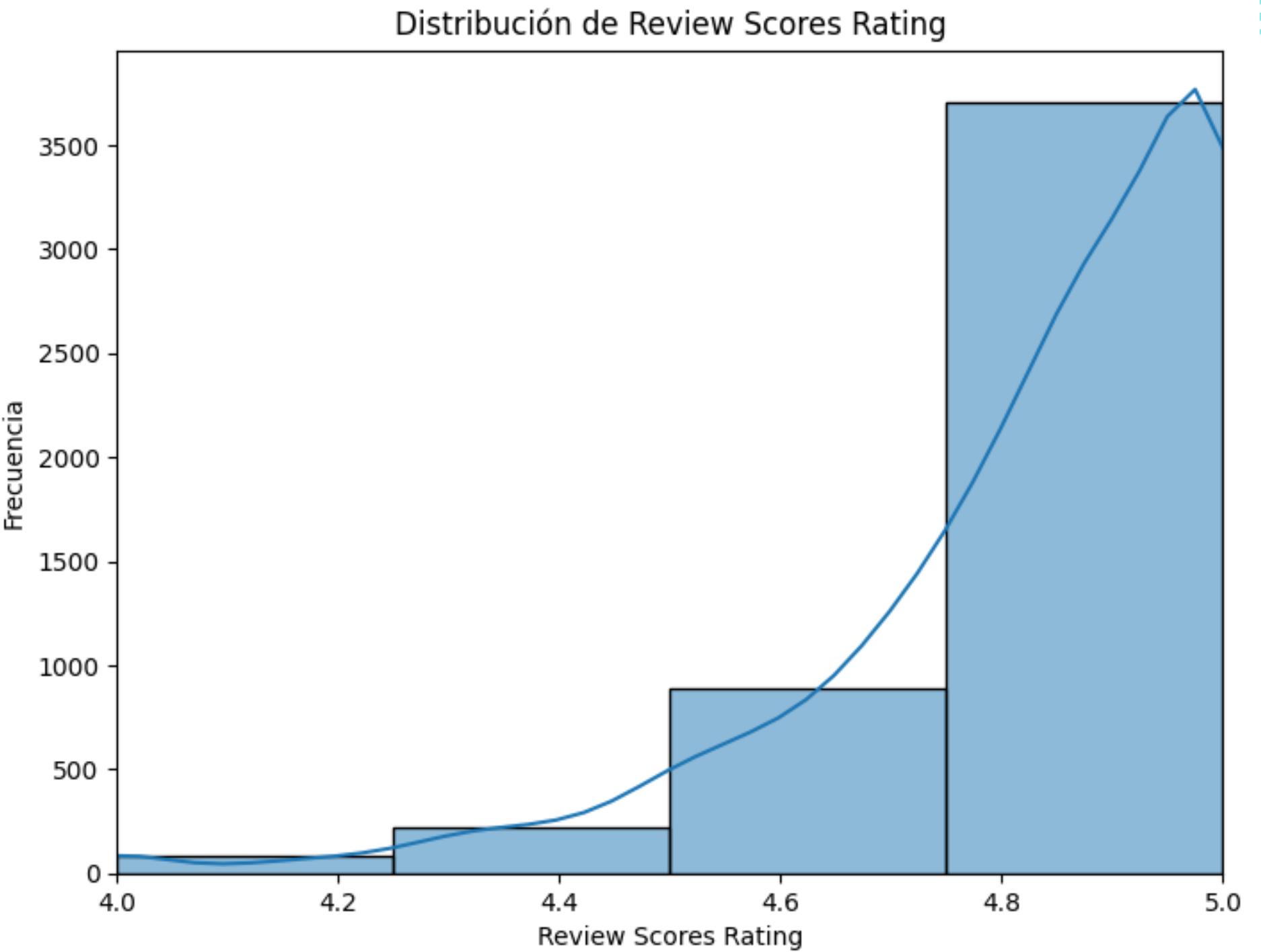
# 02. Tratamiento de Nulos y Outliers



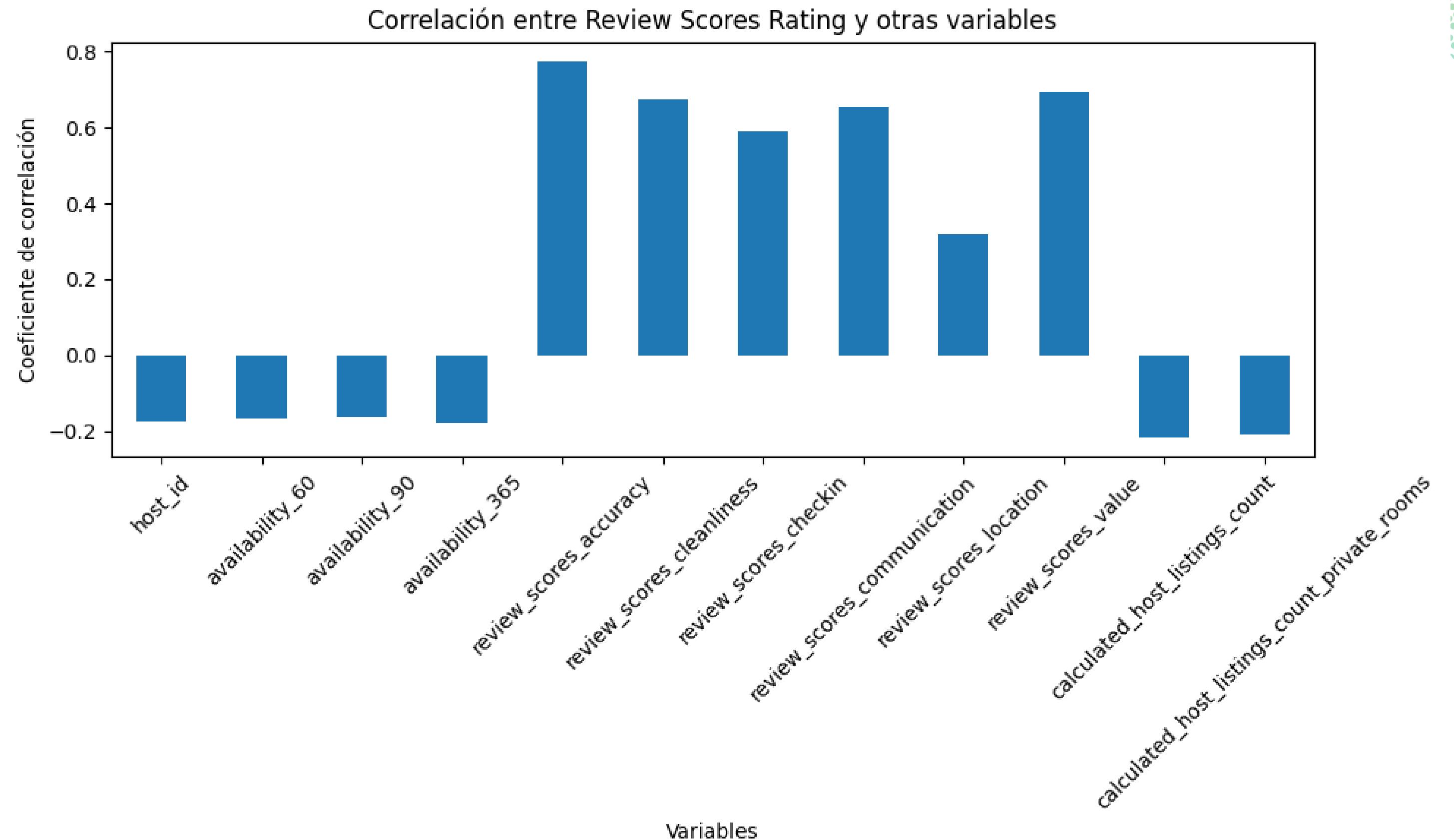
- Para imputar los valores nulos se utilizó el método de imputación por k-vecinos más cercanos
- En cuanto a los outliers no se encontraron valores que carezcan de sentido dentro de las distintas variables por lo que no fueron modificados

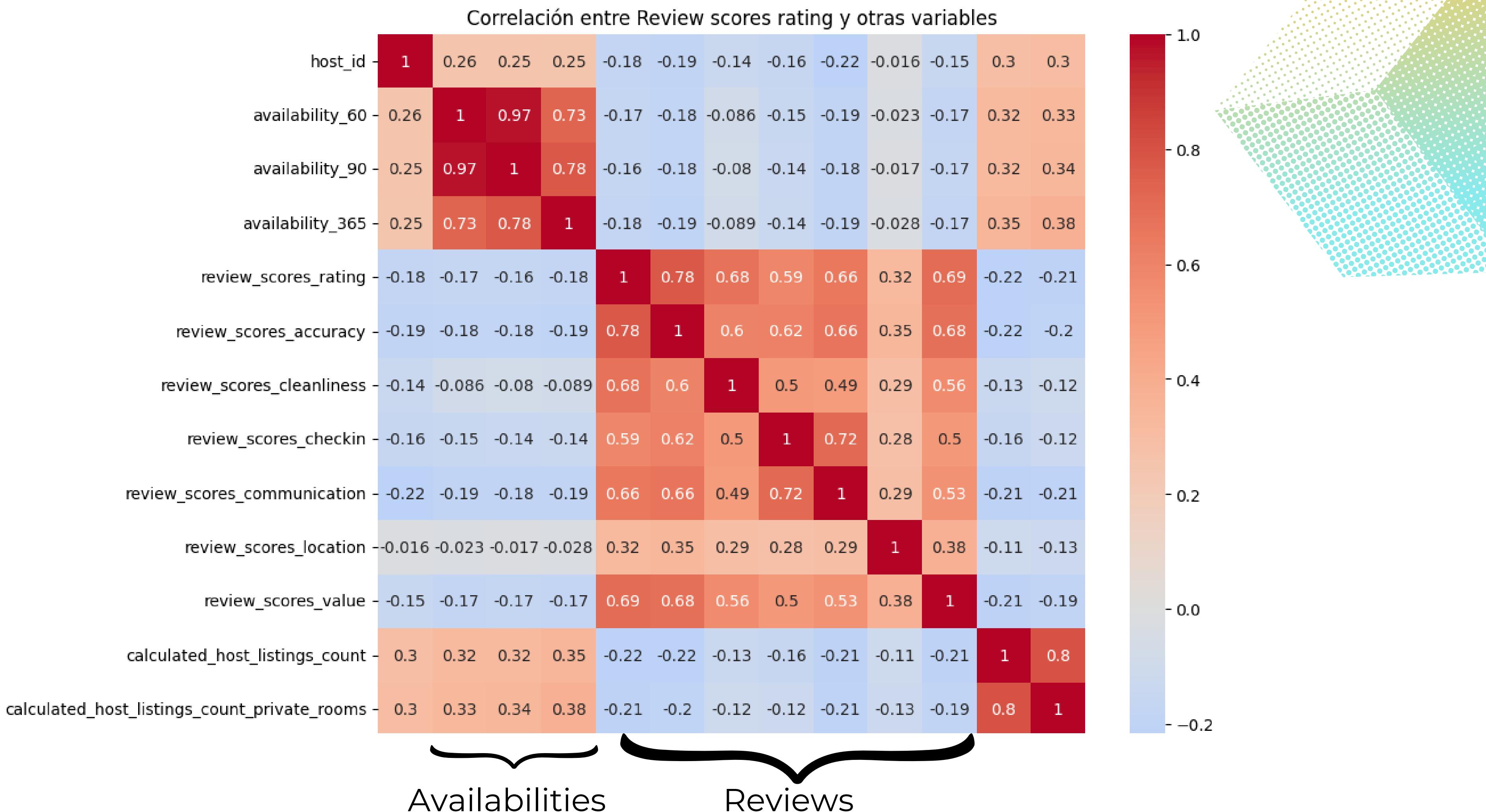
# 03. Variable objetivo

|      |          |
|------|----------|
| mean | 4.812330 |
| std  | 0.264727 |
| min  | 0.000000 |
| 25%  | 4.750000 |
| 50%  | 4.880000 |
| 75%  | 4.990000 |
| max  | 5.000000 |



# 03. Variable objetivo





# 04. Creación de variables

**Se crean nuevas variables binarias para:**

**1. Las variables true o false donde t=1 y f=0:**

host\_is\_superhost\_b, host\_has\_profile\_pic\_b, host\_identity\_verified\_b,  
has\_availability\_b, instant\_bookable\_b

**2. Aquellas que tienen dos únicos posibles valores::**

- calendar\_last\_scraped\_b donde '2022-12-05': 1, '2022-12-17': 0
- source\_b donde 'city scrape': 1, 'previous scrape': 0

**3. Las variables que estan compuestas por cadenas de texto largas si es null=0 o no=1:**

has\_description\_b, has\_neighborhood\_b,  
has\_neighbourhood\_overview\_b, has\_host\_neighbourhood\_b,  
has\_host\_about\_b

```
▶ t_f_mapping = {  
    't': 1, 'f': 0  
}  
train['host_is_superhost_b'] = train['host_is_superhost'].replace(t_f_mapping)  
train['host_has_profile_pic_b'] = train['host_has_profile_pic'].replace(t_f_mapping)  
train['host_identity_verified_b'] = train['host_identity_verified'].replace(t_f_mapping)  
train['has_availability_b'] = train['has_availability'].replace(t_f_mapping)  
train['instant_bookable_b'] = train['instant_bookable'].replace(t_f_mapping)
```

```
[ ] train['source_b'] = train['source'].map({'city scrape': 1, 'previous scrape': 0})  
train['calendar_last_scraped_b'] = train['calendar_last_scraped'].map({'2022-12-05': 1, '2022-12-17': 0})
```

# 04. Modificación de variables

```
▶ train['price'] = train['price'].str.replace('$', '')  
train['price'] = train['price'].str.replace(',', '')  
train['price'] = pd.to_numeric(train['price'])
```

```
[18] train['host_acceptance_rate'] = train['host_acceptance_rate'].str.rstrip('%').astype(float) / 100  
     train['host_response_rate'] = train['host_response_rate'].str.rstrip('%').astype(float) / 100
```

# 04. Creación de variables

**Se crearon tambien variables como:**

host\_since\_year  
cant\_amenities  
antiguedad  
host\_listings\_ratio  
host\_response\_ratio  
log\_price  
minimum\_maximum\_nights\_ratio  
description\_sentiment

# 04. Creación de variables

Y para las variables que tenían varios valores posibles se realizó frequency encoding:

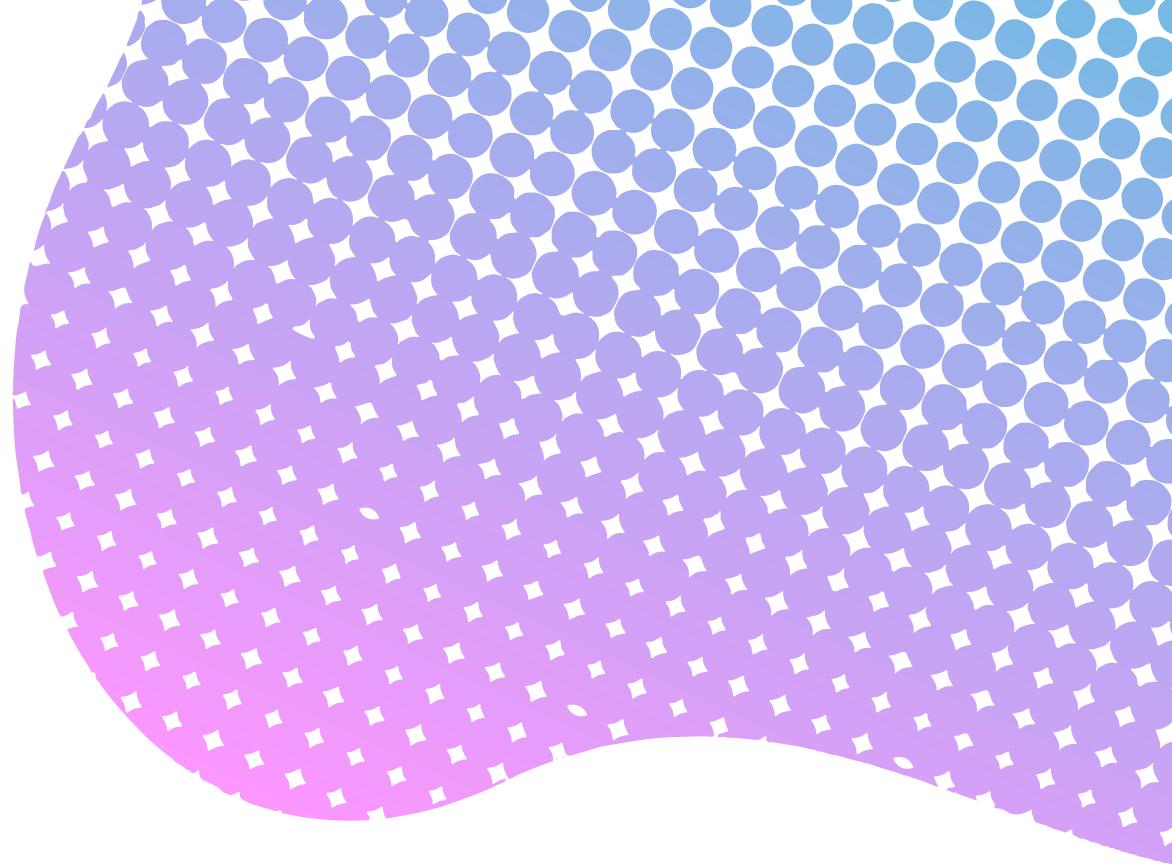
- host\_location
- host\_neighbourhood
- neighbourhood
- neighbourhood\_cleansed
- property\_type
- room\_type
- bathrooms\_text
- host\_city\_or\_state
- host\_country

```
[ ] train['room_type_num'] = train['room_type'].map({'Private room': 1, 'Entire home/apt': 2, 'Hotel room': 3, 'Shared room': 4})  
train['bathrooms_text_num'] = train['bathrooms_text'].map({'shared baths': 1, 'shared bath': 1, 'baths': 2, 'bath': 2, 'private bath': 3})
```

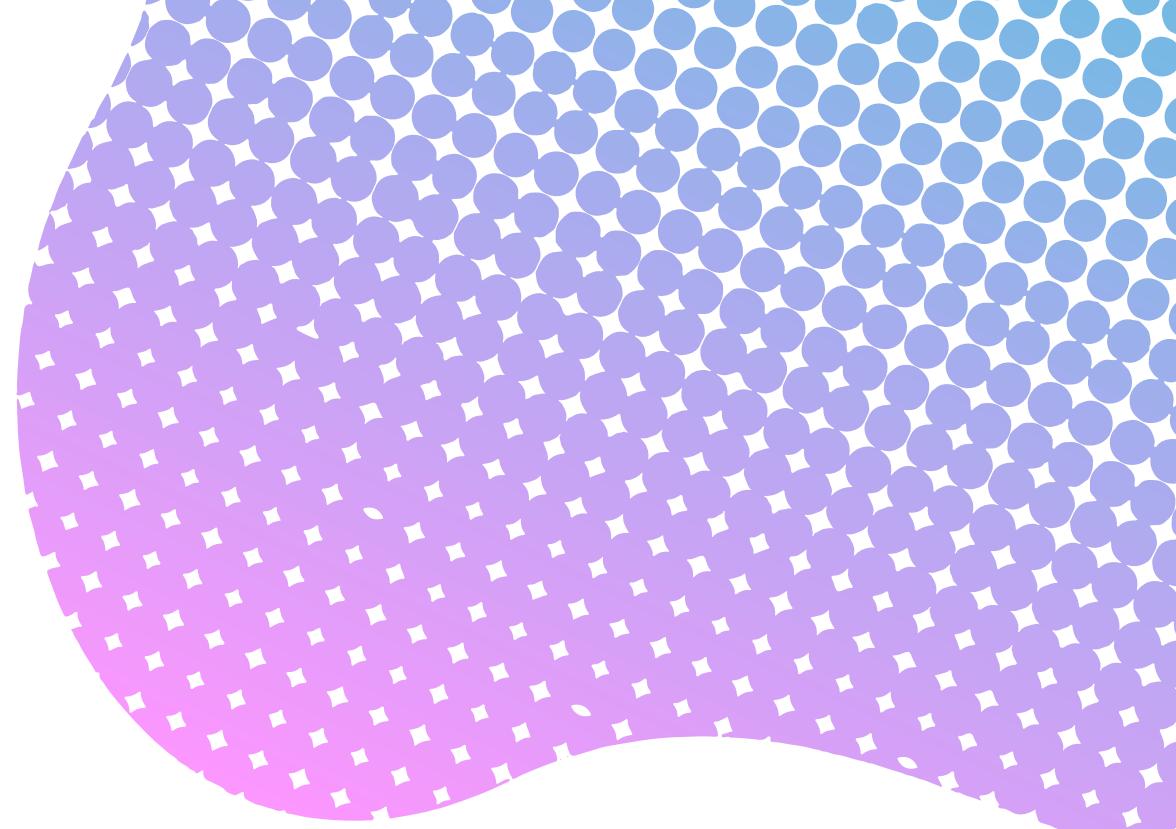
```
[ ] from sklearn.preprocessing import LabelEncoder  
neighbourhood_encoder = LabelEncoder()  
train['neighbourhood_cleaned_encoded'] = neighbourhood_encoder.fit_transform(train['neighbourhood_cleaned'])
```

# 05. Modelos y predicciones

- Se utilizaron como tamaño de las particiones a 0.15 y 0.2, ya que se encontro que eran los mejores
- Se utilizó el metodo GridSearchCV de la libreria sklearn.model\_selection como método de ajuste
  - Modelos utilizados: XGBoost, Random Forest y Regresion Ridge



# 05. Modelos y predicciones



## Random Forest

Puntajes de validación cruzada: [0.76238567 0.74664083 0.75283926 0.7820809 0.52877466]

Puntuación promedio de la validación cruzada: 0.714544264652457

{'max\_depth': None, 'min\_samples\_split': 10, 'n\_estimators': 100}



# 05. Modelos y predicciones

## XBoost

Puntajes de validación cruzada: [0.78971143 0.76931554 0.6581036 0.730894 0.77296159]  
Score: 0.7441972329749644



# 05. Modelos y predicciones

## Regresion Ridge

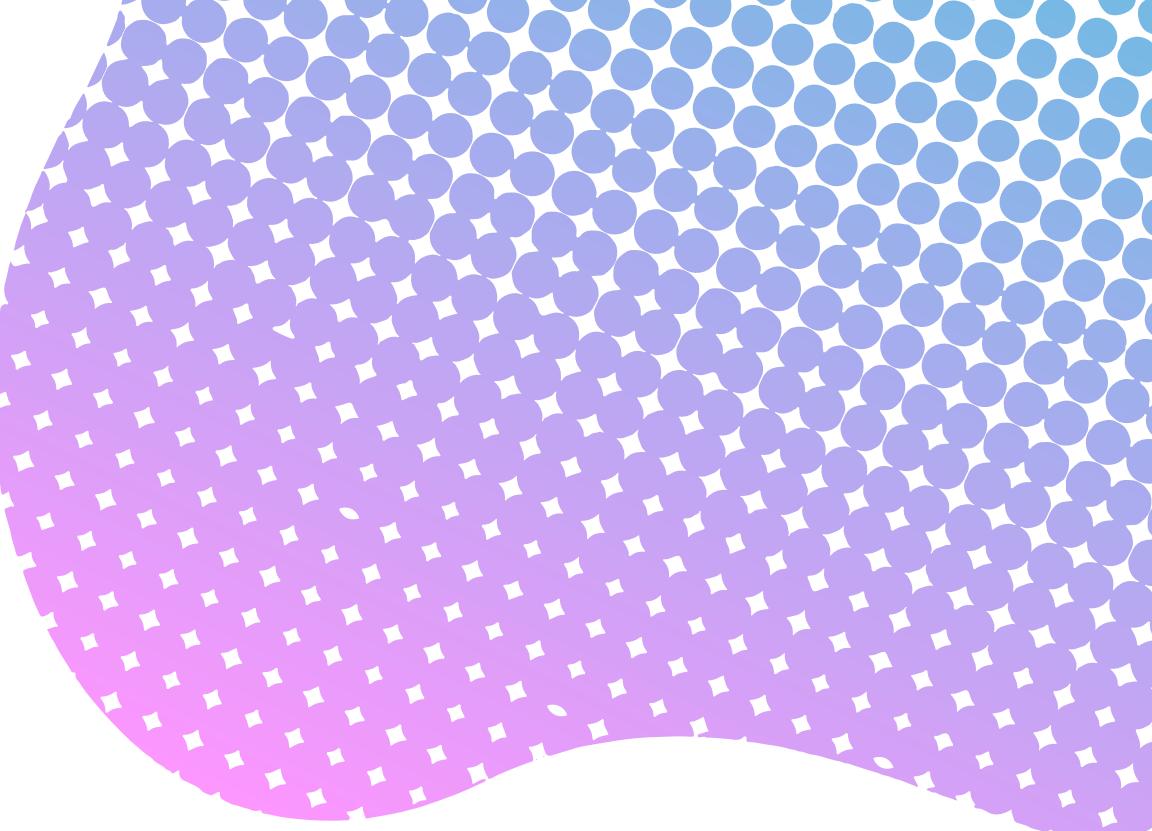


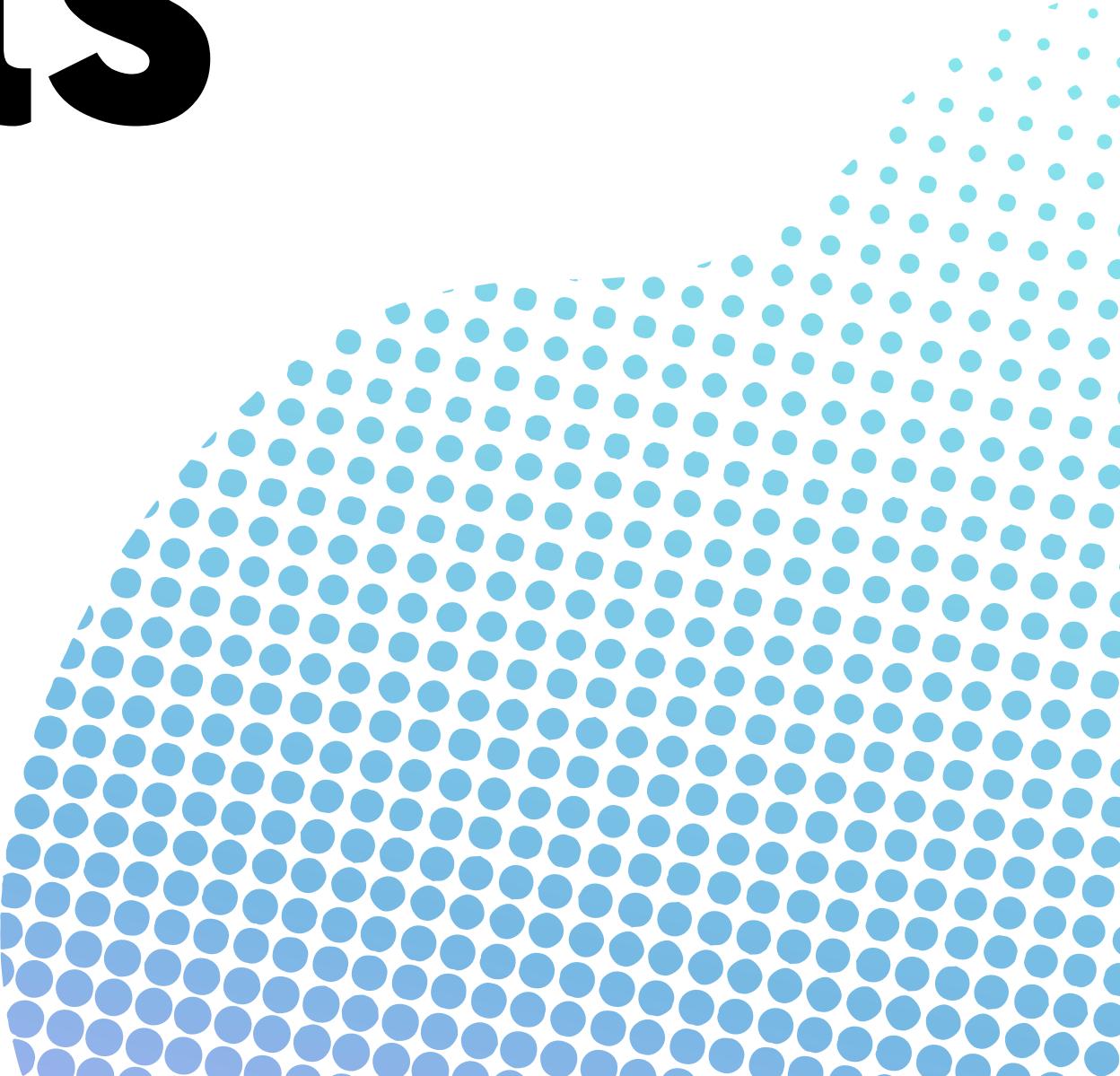
RegRidgesubmission.csv

Complete · 13d ago

0.72846

0.69743





# Gracias