

### Actividad 1: Regresión Lineal Simple y Múltiple

Para empezar trabajé con los datos nulos y outliers para que la base estuviera completa y limpia, lista para el análisis.

4)

Host acceptance rate vs host response rate

Entire home	0.53
Private room	0.54
Hotel	0.08
Shared room	0.43

Host acceptance rate vs price

Entire home	0.05
Private room	0.06
Hotel	0.04
Shared room	0.10

Host acceptance rate vs number of reviews

Entire home	0.19
Private room	0.12
Hotel	0.20
Shared room	0.15

Reviews score rating vs calculated host listings count

Entire home	0.04
Private room	0.02
Hotel	0.5

Shared room	0.16
-------------	------

Availability 365 vs number of reviews

Entire home	0.03
Private room	0.06
Hotel	0.12
Shared room	0.27

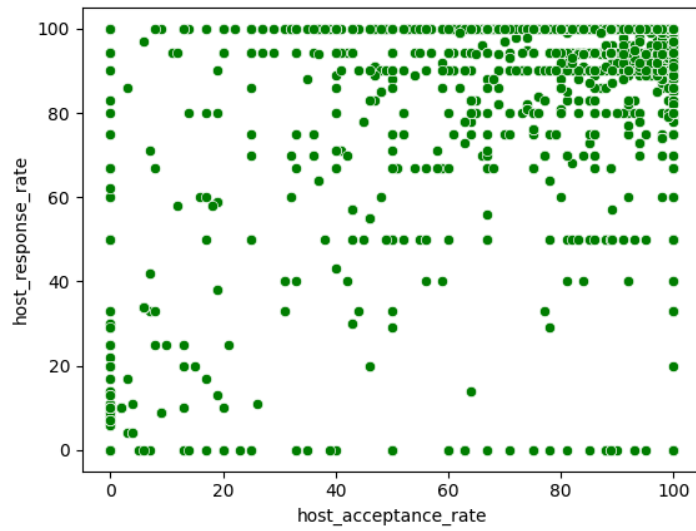
Reviews per month vs review scores communication

Entire home	0.96
Private room	0.97
Hotel	0.99
Shared room	0.97

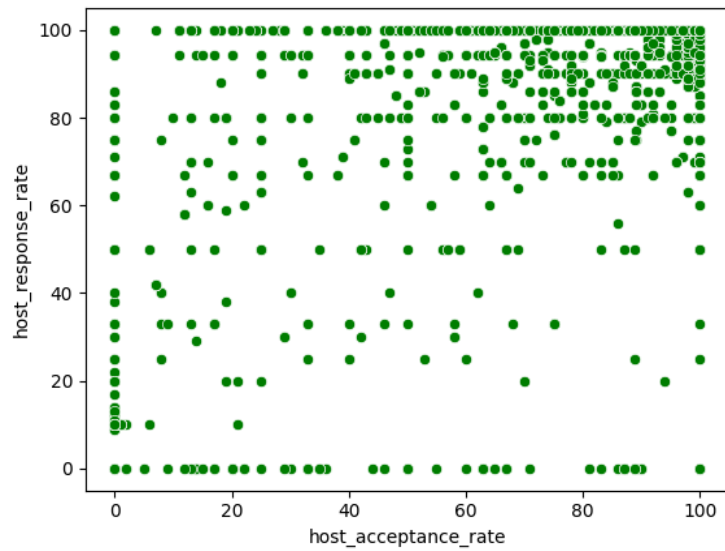
Decidí enfocar el análisis en availability 365 vs number of reviews ya que existe una gran variación que va desde 0.03 para entire home hasta 0.27 pta shared room.

En estos se puede apreciar que la tasa de aceptación se concentra en los porcentajes más elevados, aunque en entire home y private room hay más variación en general. En el caso de hotel y shared room la mayoría tienen buen host response y acceptance rate. Es importante considerar que los 0 son imputaciones hechas para los valores nulos.

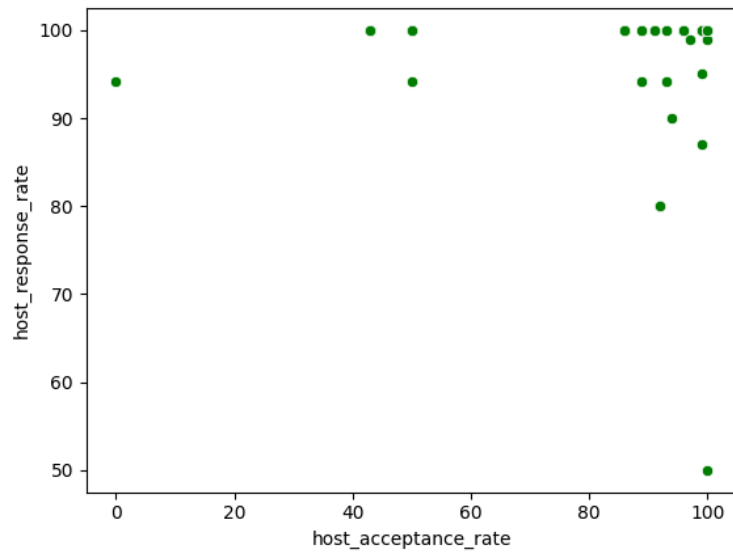
Entire home:



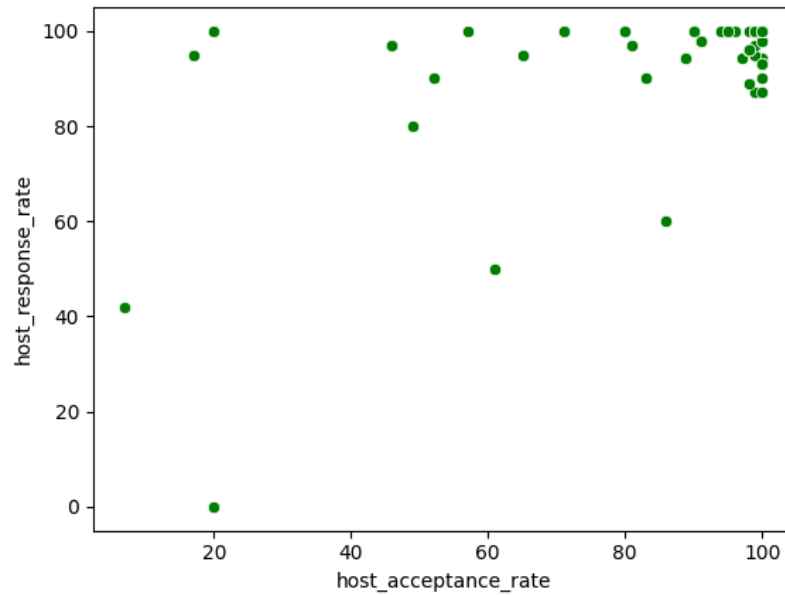
Private room:



Hotel:



Shared room:

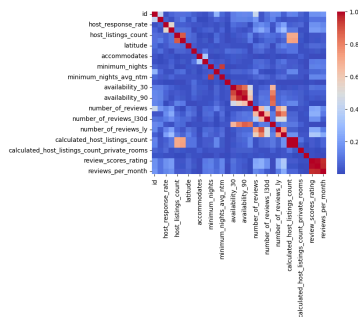


5)

Entire home:

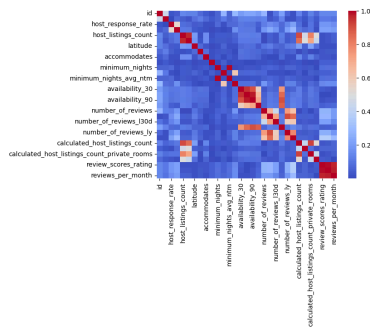
Calculated host listings count vs calculated host listings count entire homes	0.99
Review scores rating vs review score communication	0.98
Availability 60 vs availability 90	0.97

Review scores communication vs reviews per month	0.96
Review scores rating vs reviews per month	0.96
Minimum nights vs minimum nights avg ntm	0.93
Availability 30 vs availability 60	0.92
Number of reviews ltm vs number of reviews ly	0.90
Host listings count vs host total listings count	0.89
Availability 90 vs availability eoy	0.87



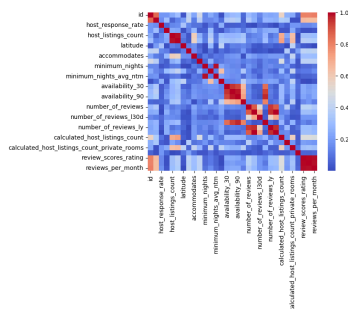
## Private room

Minimum nights vs ,minimum nights atm	0.99
Review scores rating vs review scores communication	0.99
Availability 60 vs availability 90	0.97
Review scores communication vs reviews per month	0.97
Review scores rating vs reviews per month	0.96
Host listings count vs host total listings count	0.96
Availability 30 vs availability 60	0.94
Calculated host listings count vs calculated host listings count private rooms	0.92
Number of reviews ltm vs number of reviews ly	0.91
Host listings count vs calculated host listings count	0.90



Hotel:

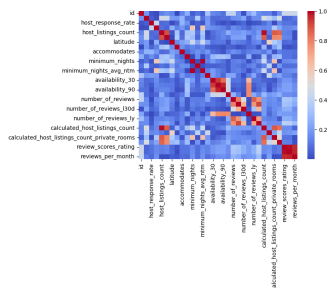
Number of reviews ltm vs estimated occupancy 365	1
Review scores rating vs review scores communication	0.99
Host listings count vs host total listings count	0.98
Minimum nights vs minimum nights avg ntm	0.98
Review scores communication vs reviews per month	0.98
Availability 60 vs availability 90	0.98
Availability 90vs availability eoy	0.96
Number of reviews ltm vs number of reviews ly	0.96
Number of reviews ly vs estimated occupancy 1365	0.96



Shared room

Minimum nights vs minimum nights avg ntm	0.99
--	------

Host listings count vs calculated host listings count	0.98
Review scores rating vs review scores communication	0.98
Availability 60 vs availability 90	0.97
Review scores communication vs reviews per month	0.96
Review scores rating vs reviews per month	0.95
Availability 30 vs availability 60	0.93
Number of reviews ltm vs estimated occupancy 1365	0.92
Host listings count vs host total listings count	0.90
Calculated host listings count vs calculated host listings count pr	0.86



De este paso puedo concluir que las variables con mayor correlación son similares en los 4 tipos de alojamientos. Algunas que me parecieron interesantes fueron review scores rating vs review scores communication ya que muestra la importancia de la atención a l@s huéspedes en la calificación final, review scores communication vs reviews per month ya que nuevamente se muestra que cuando hay buen servicio se realizan más reservas y review scores rating vs reviews per month ya que indica que mientras más reservas el anfitrión mejora y tiene mejores reseñas.

6)

Review scores rating

```

{ 'fit_intercept': True,
  'copy_X': True,
  'n_jobs': None,
  'positive': False,
  'feature_names_in_': array(['review_scores_communication', 'reviews_per_month'], dtype=object),
  'n_features_in_': 2,
  'coef_': array([ 0.89588296, -0.00868382]),
  'rank_': 2,
  'singular_': array([2646.58609394, 65.66381535]),
  'intercept_': np.float64(0.4422410436821891)}

model1.score(Vars_Indep1, Var_Dep1)

0.9820926209894811

coef_Deter1=model.score(X=Vars_Indep1, y=Var_Dep1)
coef_Deter1

0.9820926209894811

coef_Correl1=np.sqrt(coef_Deter1)
coef_Correl1

np.float64(0.9910058632467726)

```

Modelo:  $y = 0.89x_1 - 0.0006x_2 + 0.44$

Este modelo tiene un coeficiente de determinación de 0.98 y de correlación de 0.99, por lo que es muy confiable para predecir review scores rating, es decir la calificación más relevante y que resume el desempeño del anfitrión.

### Host acceptance rate

```

{ 'fit_intercept': True,
  'copy_X': True,
  'n_jobs': None,
  'positive': False,
  'feature_names_in_': array(['host_response_rate', 'availability_90',
                             'review_scores_communication'], dtype=object),
  'n_features_in_': 3,
  'coef_': array([ 0.66916452, -0.01390596, 1.77833162]),
  'rank_': 3,
  'singular_': array([5203.56572313, 2874.91863689, 262.45453054]),
  'intercept_': np.float64(19.152912020667415)}

model2.score(Vars_Indep2, Var_Dep2)

0.309984839345274

coef_Deter2=model2.score(X=Vars_Indep2, y=Var_Dep2)
coef_Deter2

0.309984839345274

coef_Correl2=np.sqrt(coef_Deter2)
coef_Correl2

np.float64(0.5567628214466857)

```

Modelo:  $y = 0.66x_1 - 0.013x_2 + 1.77x_3$

Este modelo no es muy preciso tal como lo dicen los coeficientes de determinación y correlación, de hecho este último mejoró muy poco respecto al análisis simple hecho para cada habitación en el paso 4.

### Host total listings count



```

{'fit_intercept': True,
 'copy_X': True,
 'n_jobs': None,
 'positive': False,
 'feature_names_in_': array(['host_listings_count', 'calculated_host_listings_count'],
        dtype=object),
 'n_features_in_': 2,
 'coef_': array([1.06256304, 0.63372549]),
 'rank_': 2,
 'singular_': array([13760.20835704, 3677.98404141]),
 'intercept_': np.float64(-2.338430669392636)}

model3.score(Vars_Indep3, Var_Dep3)

0.8080123085712159

coef_Deter3=model3.score(X=Vars_Indep3, y=Var_Dep3)
coef_Deter3

0.8080123085712159

coef_Correl3=np.sqrt(coef_Deter3)
coef_Correl3

np.float64(0.8988950486965739)

```

Modelo:  $y = 1.06x_1 + 0.63x_2 - 2.34$

Considero que esta es una variable difícil de predecir ya que aunque los coeficientes son altos, estos es debido a que estoy usando variables que habían tenido alta correlación pero que realmente son muy similares y no aportan valor, estas son host listings count y calculated host listings count.

Referencia del dataset:

*Get the Data.* (s.f.). Inside Airbnb. <https://insideairbnb.com/get-the-data/>