

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

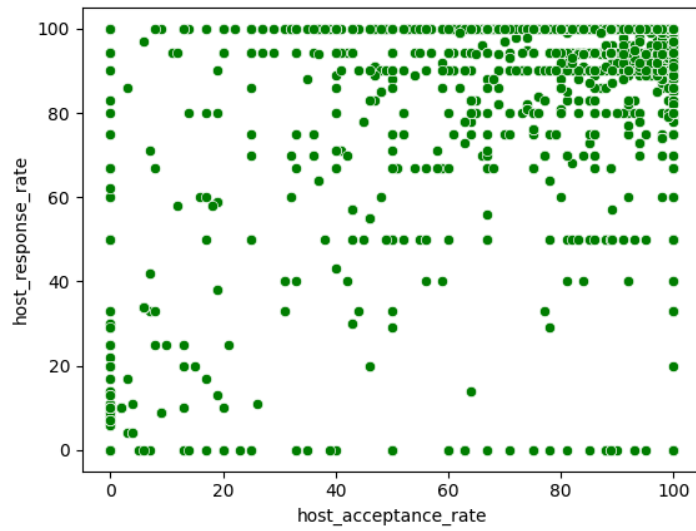
Reviews per month vs review scores communication

Entire home	0.96
Private room	0.97
Hotel	0.99
Shared room	0.97

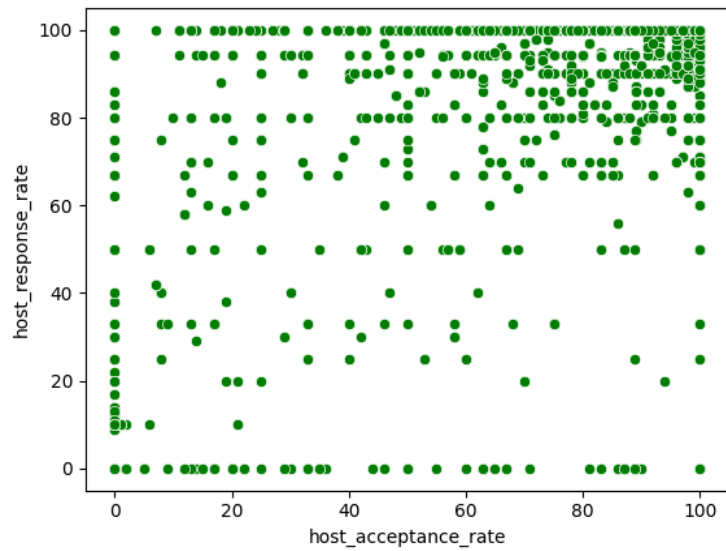
Decidí enfocar el análisis en availability 365 vs number of reviews ya que existe una gran variación que va desde 0.03 para entire home hasta 0.27 para shared room.

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.

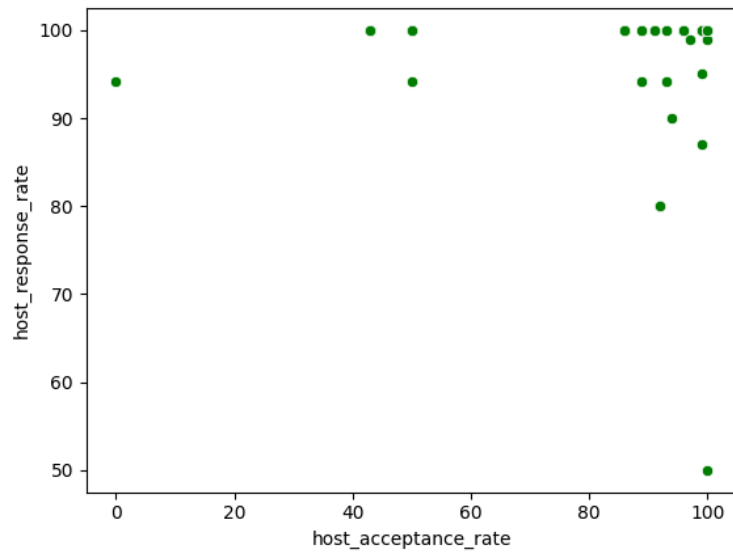
Entire home:



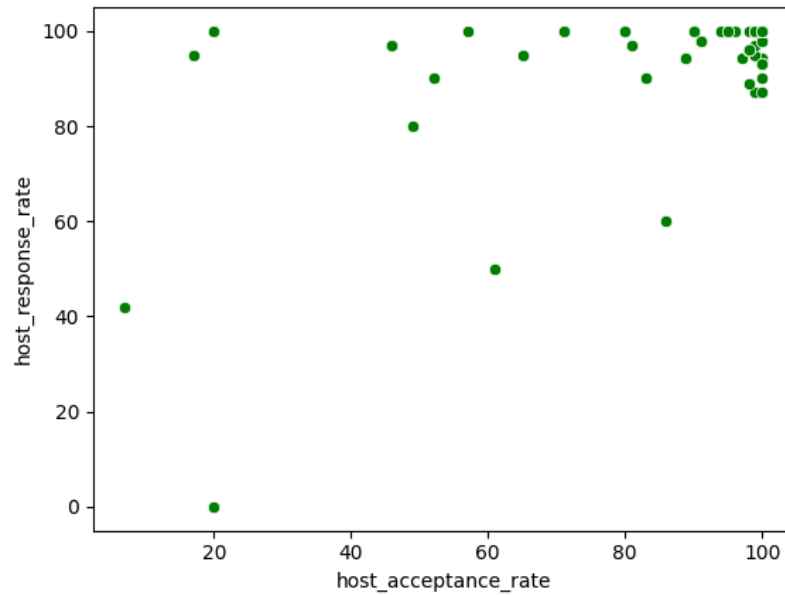
Private room:



Hotel:



Shared room:

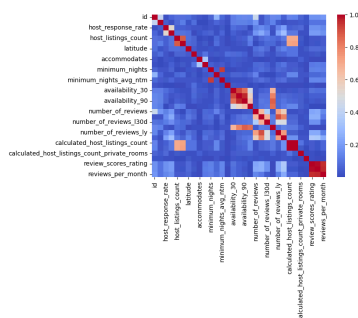


5)

Entire home:

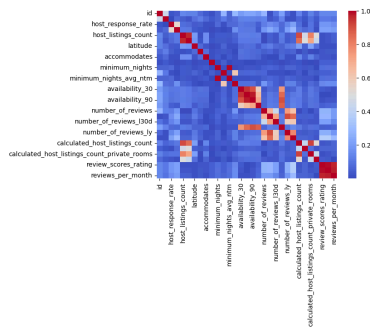
Calculated host listings count vs calculated host listings count entire homes	0.99
Review scores rating vs review score communication	0.98
Availability 60 vs availability 90	0.97

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



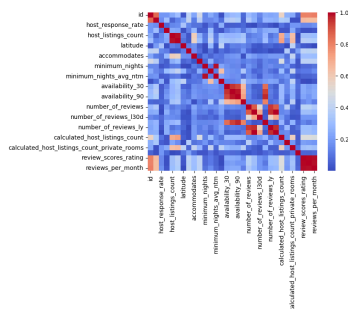
Private room

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:

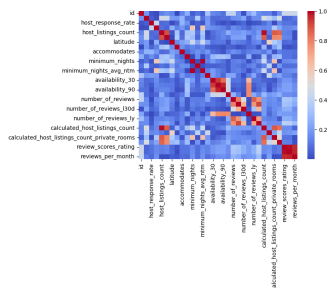
Number of reviews ltm vs estimated occupancy 365	1
Review scores rating vs review scores communication	0.99
Host listings count vs host total listings count	0.98
Minimum nights vs minimum nights avg ntm	0.98
Review scores communication vs reviews per month	0.98
Availability 60 vs availability 90	0.98
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 l365	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)

Nota: para identificar las correlaciones de las regresiones simples estuve usando la información de entire home debido a que es el tipo de cuarto más común, sin embargo, varía en cada uno de ellos, lo que afecta el resultado final.

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: $y = 0.66x_1 - 0.013x_2 + 1.77x_3 + 19.15$

Predecir la tasa de aceptación que tienen los anfitriones es difícil debido a las bajas correlaciones que existen con esta variable, sin embargo, al combinar host acceptance rate, availability 90 y review scores communication se alcanza 0.3 en determinación y .56 en correlación.

Host is superhost


```

{'fit_intercept': True,
 'copy_X': True,
 'n_jobs': None,
 'positive': False,
 'feature_names_in_': array(['estimated_occupancy_l365d', 'host_acceptance_rate',
                             'host_response_rate'], dtype=object),
 'n_features_in_': 3,
 'coef_': array([0.00158156, 0.00254775, 0.00177148]),
 'rank_': 3,
 'singular_': array([14826.52662746, 3965.37347975, 2155.98276612]),
 'intercept_': np.float64(-0.16037021078736952)}

model3.score(Vars_Indep3, Var_Dep3)

0.14814258333213637

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

0.14814258333213637

coef_Correl3=np.sqrt(coef_Deter3)
coef_Correl3

np.float64(0.3848929504838149)

```

Modelo: $y = 0.001x_1 + 0.003x_2 + 0.002x_3 - 0.16$

La variable si host es superhost no cuenta con muy altas correlaciones, sin embargo, al aplicar las más altas (estimated occupancy l365d, host acceptance rate y host response rate) se puede subir la determinación a 0.148 y la correlación a 0.38.

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

Esta variable ya contaba con coeficientes altos, que todavía aumentan gracias al modelo múltiple con las variables independientes host listings count y calculated host listings count, la determinación y correlación son 0.88 y 0.89 respectivamente.

Accommodates

```
{'fit_intercept': True,
 'copy_X': True,
 'n_jobs': None,
 'positive': False,
 'feature_names_in_': array(['bedrooms', 'estimated_revenue_l365d', 'price'], dtype=object),
 'n_features_in_': 3,
 'coef_': array([-2.46073246e-02,  2.60695523e-06,  3.11415937e-06]),
 'rank_': 3,
 'singular_': array([4.04144248e+07,  2.79614505e+06,  1.41638912e+03]),
 'intercept_': np.float64(2.9484937459508633)}
```

```
model5.score(Vars_Indep5, Var_Dep5)
```

```
0.07493534599919083
```

```
coef_Deter5=model5.score(X=Vars_Indep5, y=Var_Dep5)
coef_Deter5
```

```
0.07493534599919083
```

```
coef_Correl5=np.sqrt(coef_Deter5)
coef_Correl5
```

```
np.float64(0.27374321178650407)
```

Modelo: $y = -2.46x_1 + 2.60x_2 + 3.11x_3 + 2.49$

Pocas variables están muy correlacionadas con accommoates, sin embargo escogí las más altas (bedrooms, estimated revenue l365d y price) para obtener los coeficientes de 0.07 y 0.27 respectivamente. Algo interesante es que cuando aumentan los cuartos se resta en el modelo, disminuyendo accommoates.

Bedrooms

```
{'fit_intercept': True,
 'copy_X': True,
 'n_jobs': None,
 'positive': False,
 'feature_names_in_': array(['accommodates', 'bathrooms', 'beds'], dtype=object),
 'n_features_in_': 3,
 'coef_': array([-0.16432622, -0.00633992,  0.26863878]),
 'rank_': 3,
 'singular_': array([3670.29511351,  393.21223973,  218.30533299]),
 'intercept_': np.float64(1.7092302315239043)}
```

```
model6.score(Vars_Indep6, Var_Dep6)
```

```
0.21755461625113015
```

```
coef_Deter6=model6.score(X=Vars_Indep6, y=Var_Dep6)
coef_Deter6
```

```
0.21755461625113015
```

```
coef_Correl6=np.sqrt(coef_Deter6)
coef_Correl6
```

```
np.float64(0.4664275037464345)
```

Modelo: $y = -0.16x_1 - 0.006x_2 + 0.27x_3 + 1.7$

Aunque en general la correlación con bedrooms no es muy alta, combinar accommoates, bathrooms y beds que son las más relevantes (y que además se relacionan de forma lógica) lleva a .22 de determinación y 0.47 de correlación. Algo interesante es que tanto accommoates como bathrooms restan en el modelo, cuando se podría considerar al revés.

Price

```

{'fit_intercept': True,
 'copy_X': True,
 'n_jobs': None,
 'positive': False,
 'feature_names_in_': array(['estimated_revenue_l365d', 'accommodates', 'bedrooms'],
    dtype=object),
 'n_features_in_': 3,
 'coef_': array([ 7.57034701e-03,  1.82089437e+02, -3.63665556e+01]),
 'rank_': 3,
 'singular_': array([4.04132243e+07, 1.41715706e+03, 3.65481363e+02]),
 'intercept_': np.float64(168.40355876053445)}

model7.score(Vars_Indep7, Var_Dep7)

0.013183347544848845

coef_Deter7=model7.score(X=Vars_Indep7, y=Var_Dep7)
coef_Deter7

0.013183347544848845

coef_Correl7=np.sqrt(coef_Deter7)
coef_Correl7

np.float64(0.11481875955108052)

```

Modelo: $y = 7.54x_1 + 1.82x_2 - 3.64x_3 + 168.4$

A pesar de haber combinado las variables más correlacionadas, los coeficientes para price siguen siendo muy bajos (0.01 y 0.11 respectivamente). Considero que esto se debe a que en otros tipos de cuarto hay variables más relevantes que estas, es decir que el dataset varía considerablemente.

Review scores value

```

{'fit_intercept': True,
 'copy_X': True,
 'n_jobs': None,
 'positive': False,
 'feature_names_in_': array(['review_scores_rating', 'review_scores_accuracy',
    'review_scores_communication'], dtype=object),
 'n_features_in_': 3,
 'coef_': array([0.58804088, 0.30876592, 0.09217814]),
 'rank_': 3,
 'singular_': array([458.80563931, 28.00205405, 20.25407839]),
 'intercept_': np.float64(-0.0005005071308294973)}

model8.score(Vars_Indep8, Var_Dep8)

0.9904883860114281

coef_Deter8=model8.score(X=Vars_Indep8, y=Var_Dep8)
coef_Deter8

0.9904883860114281

coef_Correl8=np.sqrt(coef_Deter8)
coef_Correl8

np.float64(0.995232830051053)

```

Modelo: $y = 0.59x_1 + 0.31x_2 + 0.09x_3 - 0.0005$

La variable de review scores value se relaciona fuertemente con otras reviews, al combinarla se pueden alcanzar coeficientes casi perfectos de 0.99 tanto para determinación como para correlación.

Bathrooms

```

{'fit_intercept': True,
 'copy_X': True,
 'n_jobs': None,
 'positive': False,
 'feature_names_in_': array(['beds', 'availability_90', 'availability_365'], dtype=object),
 'n_features_in_': 3,
 'coef_': array([ 9.92766646e-01, -3.58842462e-04, -8.79541810e-04]),
 'rank_': 3,
 'singular_': array([20532.39170053, 3879.21275364, 2222.82579997]),
 'intercept_': np.float64(-0.18460978820628426)}

model9.score(Vars_Indep9, Var_Dep9)

0.9830014211215022

coef_Deter9=model9.score(X=Vars_Indep9, y=Var_Dep9)
coef_Deter9

0.9830014211215022

coef_Correl9=np.sqrt(coef_Deter9)
coef_Correl9

np.float64(0.9914642813140079)

```

Modelo: $y = 9.93x_1 - 3.59x_2 - 8.79x_3 - 0.18$

El modelo múltiple para bathrooms también es muy bueno, ya que al combinar variables que tienen alta correlación (beds, availability 90 y availability 365) se obtienen coeficientes de 0.98 para determinación y 0.99 para correlación. Las variables de availability me llaman la atención ya que cuando un alojamiento tiene más disponibilidad disminuye el número de baños con los que cuenta.

Reviews per month

```

{'fit_intercept': True,
 'copy_X': True,
 'n_jobs': None,
 'positive': False,
 'feature_names_in_': array(['review_scores_checkin', 'review_scores_communication',
 'review_scores_location'], dtype=object),
 'n_features_in_': 3,
 'coef_': array([-2.77654914, -0.78885786, -6.12083891]),
 'rank_': 3,
 'singular_': array([460.58391056, 33.28280758, 23.79749508]),
 'intercept_': np.float64(48.83672540905934)}

model10.score(Vars_Indep10, Var_Dep10)

0.9562957324386572

coef_Deter10=model10.score(X=Vars_Indep10, y=Var_Dep10)
coef_Deter10

0.9562957324386572

coef_Correl10=np.sqrt(coef_Deter10)
coef_Correl10

np.float64(0.9779037439536967)

```

Modelo: $y = -2.78x_1 - 0.79x_2 - 6.12x_3 + 48.84$

Reviews per month está muy correlacionada con las distintas reviews de forma negativa, es decir que cuando aumenta la calificación disminuye el número de reviews. El modelo múltiple es muy bueno, aunque mejora poco respecto a los simples, su determinación es de 0.96 y su correlación de 0.98.

Referencia del dataset:

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