

Cristian Vega Castro

Alojamientos en Barcelona

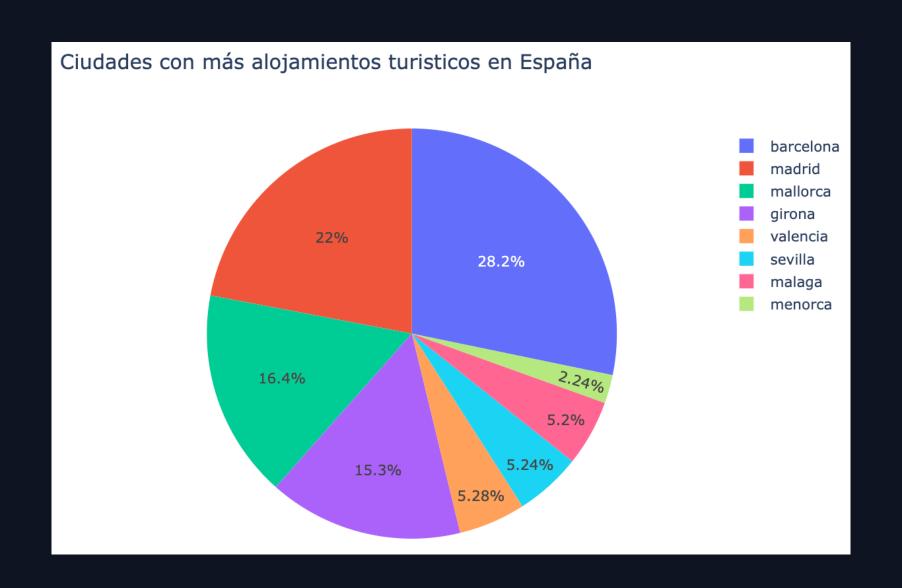
Cada año que pasa, las personas tienen más posibilidades de viajar a otros países y uno de los más atractivos para los países sudamericanos es España, principalmente por el idioma. España se caracteriza por tener ciudades muy turísticas en el verano europeo, como Barcelona, Madrid, la isla de Mallorca, entre otras. Existen muchas aplicaciones para alojar en distintas ciudades del mundo, dentro de las cuales su precio depende de muchos factores, como cercanía a atractivos turísticos, los servicios que ofrecen, la capacidad de personas admitidas, etc.

Este proyecto busca generar un modelo que permita estimar el precio de renta de alojamientos turísticos, ya sea de una habitación privada o de la renta de una casa/apartamento en la ciudad de Barcelona, comprendiendo las relaciones entre las características de los alojamientos y sus respectivas rentas, con el fin de ayudar a los anfitriones a ajustar sus valores de renta y a los que quieran incursionar en esta área de arrendar alojamientos turísticos tomando la mejor decisión.

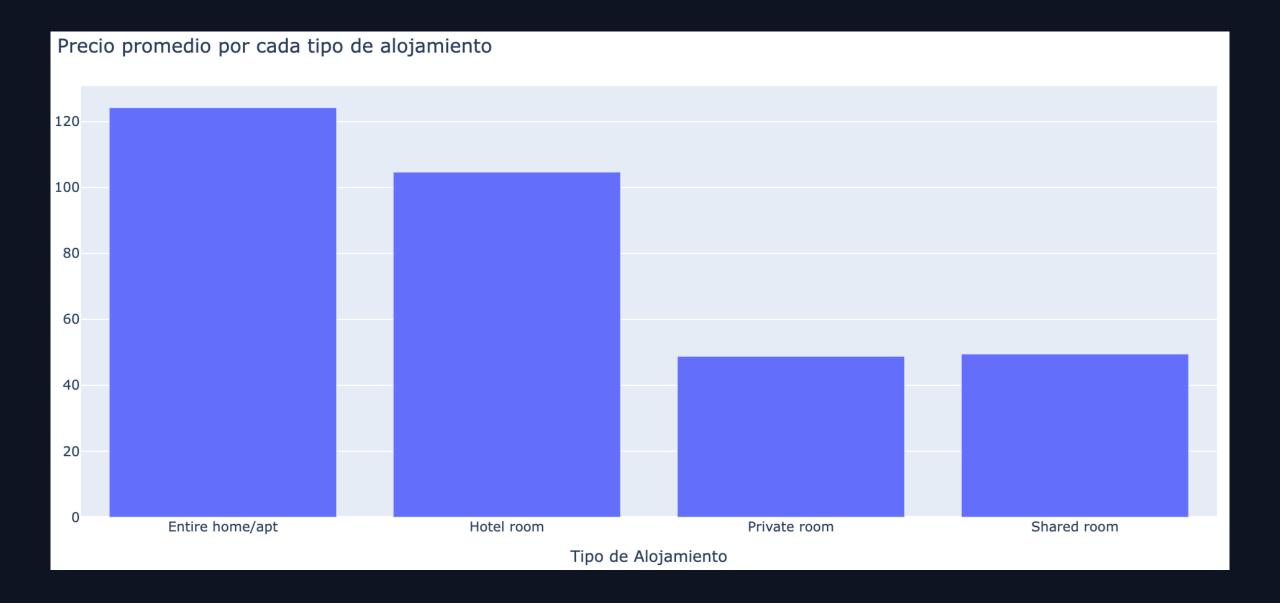




Datos iniciales



Datos iniciales



Resumen de los datos originales

Data Summary									
dataframe	Values								
Number of rows Number of columns	10000 18								

Data Types								
Column Type	Count							
int64 float64 string bool	8 6 2 2							

bool									
column_name	true	true rate	hist						
is_instant_bookable has_availability	5806 10000	0.58 1	1 1						

En primera instancia, se seleccionaron los features que se creían importantes para el target, obteniendo el siguiente resumen de cada componente del modelo:

Dado que 'has_availability' solo está compuesto por valores True, eliminamos esta variable

				number						
column_name	NA	NA %	mean	sd	р0	p25	p50	p75	p100	hist
price	254	2.54	100	86	6	45	75	120	500	
latitude	0	0	40	1.4	37	40	40	41	42	
longitude	0	0	0.28	3.2	-6	-3.7	2.2	2.9	4.3	
cant_comodidades	0	0	21	11	1	13	19	27	99	
accommodates	0	0	4.3	2.6	1	2	4	6	29	
bathrooms	74	0.74	1.6	0.99	0	1	1	2	13	
bedrooms	70	0.7	1.9	1.4	0	1	2	3	50	
beds	45	0.45	2.9	2.3	0	1	2	4	30	L
minimum_nights	0	0	5	18	1	1	2	4	1100	
maximum_nights	0	0	760	500	1	62	1100	1100	1100	
availability_30	0	0	13	12	0	0	11	26	30	
availability_60	0	0	29	23	0	2	29	53	60	
availability_90	0	0	46	34	0	10	49	80	90	
availability_365	0	0	190	130	0	66	180	320	360	

			string	
column_name	NA	NA %	words per row	total words
room_type city	0 0	0	2 1	20000 10000

Dataset acotado

Ya que Barcelona junto con Entire home/apt y Private room son los datos con mayor frecuencia en el dataset, le daremos énfasis a estos para crear un nuevo conjunto de datos, manteniendo todas las columnas anteriormente mostradas, pero eliminando la columna 'city' y encodeando 'room_type':

• 0 o False: Entire home/apt

1 o True: Private room

Data Summary								
dataframe	Values							
Number of rows Number of columns	2678 16							

Data Types								
Column Type	Count							
int64 float64 bool	8 6 2							

	number									
column_name	NA	NA %	mean	sd	p0	p25	p50	p75	p100	hist
price latitude longitude cant_comodidades accommodates bathrooms bedrooms beds minimum_nights maximum_nights availability_30 availability_60 availability_90 availability_365	43 0 0 0 66 15 17 0 0 0	1.61 0 0 0 2.46 0.56 0.63 0 0	84 41 2.2 19 3.3 1.4 1.6 2.2 7 710 12 27 44 180	76 0.015 0.018 9.9 2.3 0.64 0.95 1.9 14 510 11 23 34 140	6 41 2.1 1 0 1 1 1 0 0	35 41 2.2 12 2 1 1 1 60 0 2 8 49	60 41 2.2 17 2 1 1 2 2 1100 9 24 46 180	100 41 2.2 24 4 2 2 3 1100 23 50 78 320	500 41 2.2 99 16 8 8 20 300 1100 30 60 90 360	IIII

bool									
column_name	true	true rate	hist						
is_instant_bookable room_type	1309 1442	0.49 0.54							

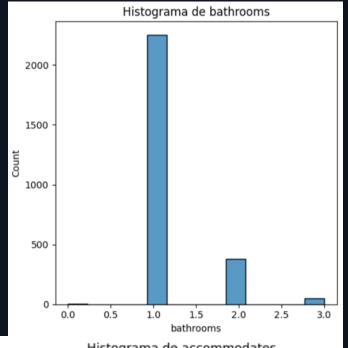
Dataset final

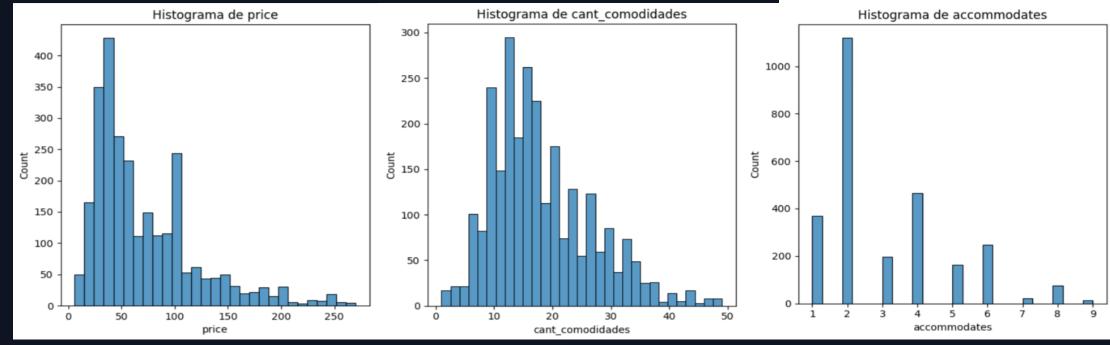
Para identificar los *outliers*, primero se dividió el dataset en 2 conjuntos de datos: Alojamientos que pertenecen a Entire home/apt y Alojamientos que pertenecen a Private room, esto con el motivo de que tienen características muy distintas y un dato común de un tipo de habitación podría ser *outlier* del otro. El método utilizado para identificar estos *outliers* fue IQR, eliminando solo este dato extremo (y así no eliminar la fila completa) y para imputar los valores faltantes se empleó la mediana ya que los datos no tenían una distribución normal, por lo que el promedio no pudo ser utilizado.

number										
column_name	NA	NA %	mean	sd	p0	p25	p50	p75	p100	hist
price latitude longitude cant_comodidades accommodates bathrooms bedrooms beds minimum_nights maximum_nights availability_30 availability_60 availability_90 availability_365	000000000000000000000000000000000000000	000000000000000000000000000000000000000	71 41 2.2 18 3.1 1.2 1.5 1.9 2.2 710 12 27 44 180	48 0.015 0.018 8.6 1.8 0.42 0.92 1.4 1.1 510 11 23 34 140	6 41 2.1 1 0 1 1 1 0 0	35 41 2.2 12 2 1 1 1 60 0 2 8 49	55 41 2.2 16 2 1 1 2 1100 9 24 46 180	99 41 2.2 23 4 1 2 3 1100 23 50 78 320	270 41 2.2 49 9 3 6 7 8 1100 30 60 90 360	4444-447333

Análisis univariado

A continuación, se muestra el análisis univariado de algunas variables del modelo con el nuevo dataset que corresponde a la ciudad de Barcelona y tipo de habitación Entire home/apt y Private room

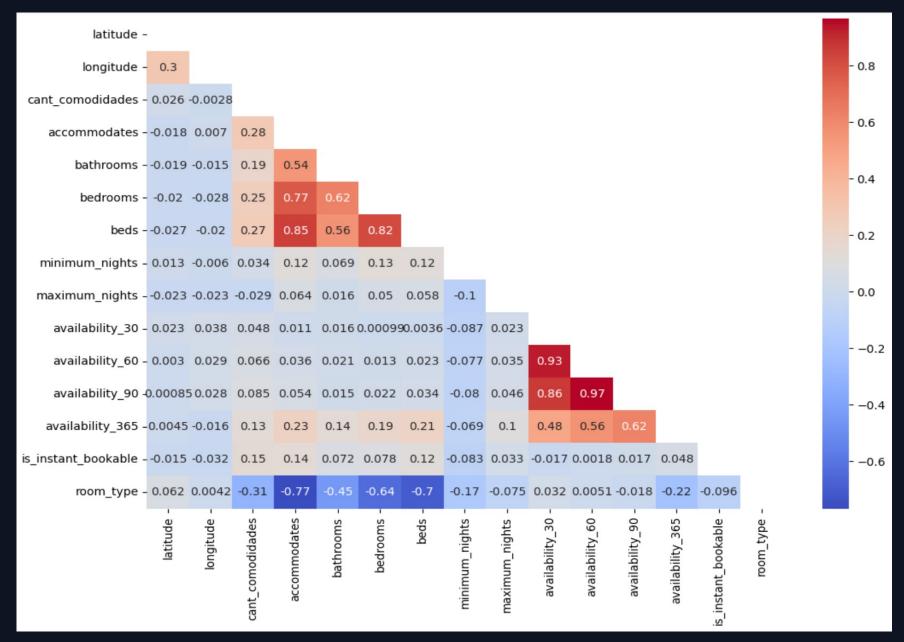




Dejando de lado el target, se aprecia una cierta correlación entre algunas variables, por lo que se quitan del modelo (beds, availability_60, y availability_90).

Además, para no generar demasiada variabilidad en el modelo, se eliminan las columnas longitude y latitude)

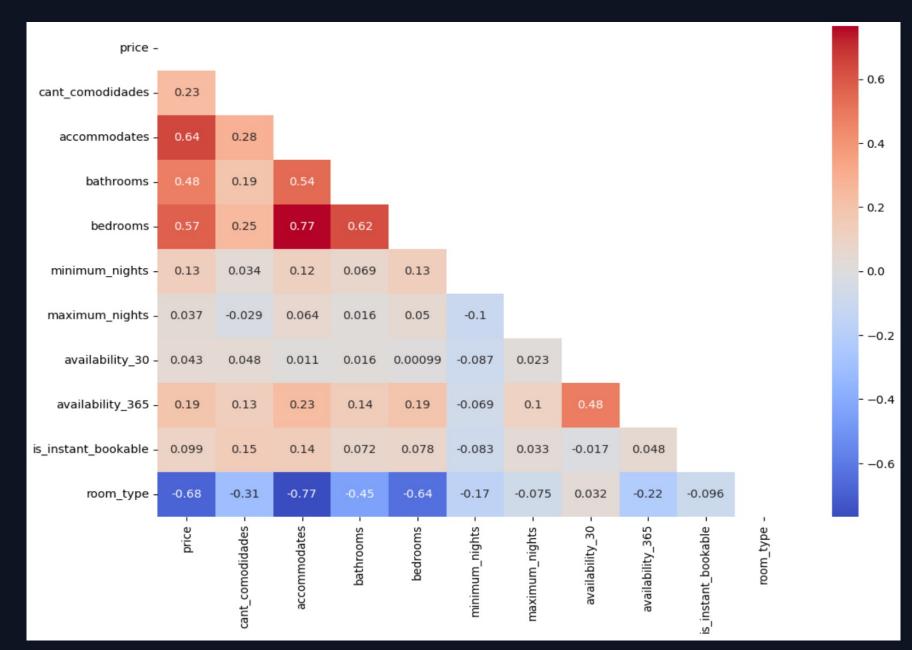
Análisis bivariado



Análisis bivariado

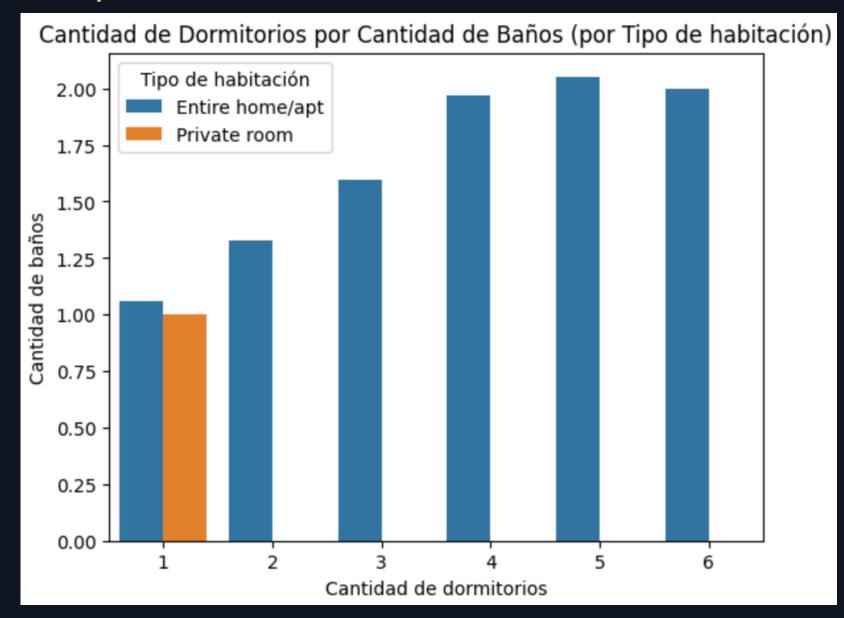
Se realiza un nuevo análisis bivariado, considerando el *target* y los features finales.

En este caso, ya no se presenta una gran correlación entre los features y además, el target price no está altamente correlacionado a ninguna variable.

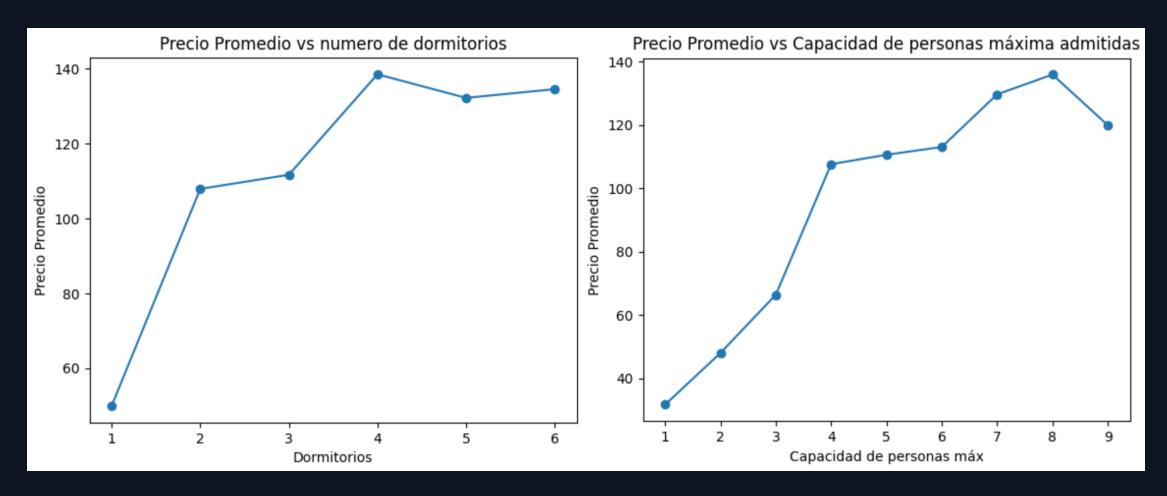


Dormitorios por cantidad de baños

Según vemos, para el caso de Private room, como su nombre lo dice, solo posee 1 dormitorio y en promedio ofrece 1 baño. En el caso de Entire home/apt, este cuando posee 1 dormitorio, también ofrece un baño, pero a medida que aumenta la cantidad de dormitorios, en promedio, aumenta la cantidad de baños hasta un máximo de 5 dormitorios



Respuestas a las hipótesis de interés



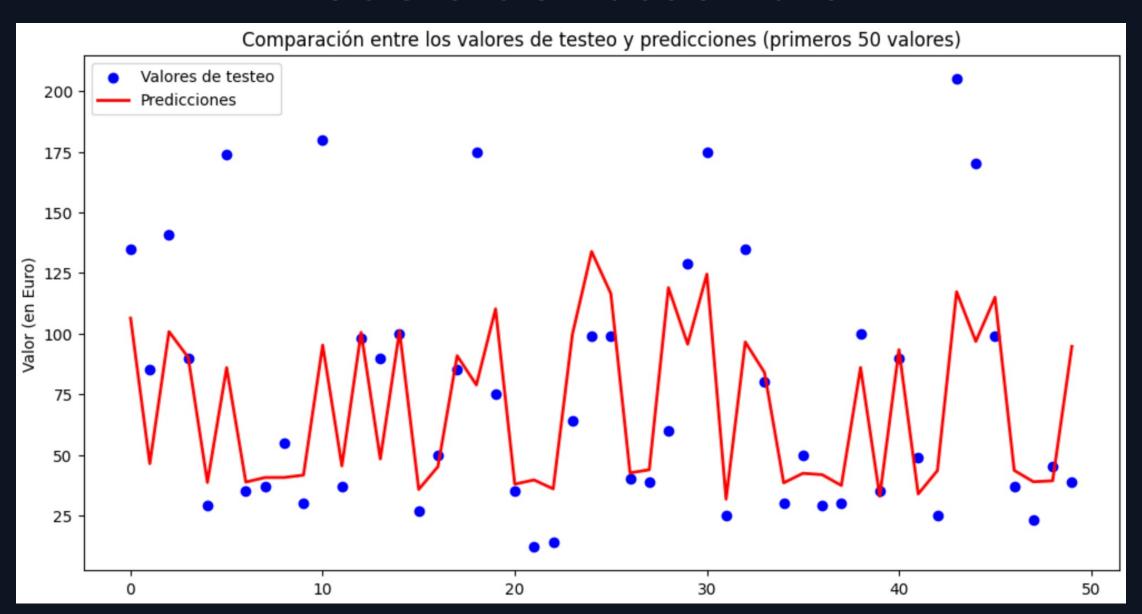
Según vemos en el gráfico de la izquierda, existe una correlación positiva entre el precio promedio de un alojamiento y la cantidad de dormitorios que este ofrece, así como en el gráfico de la derecha, que muestra que también existe una correlación positiva (no tan marcada) entre el precio promedio del alojamiento y la cantidad de personas que permite, es decir, a mayor precio del alojamiento, en general, este tendrá más dormitorios y aceptará más personas.

Modelos de Machine Learning

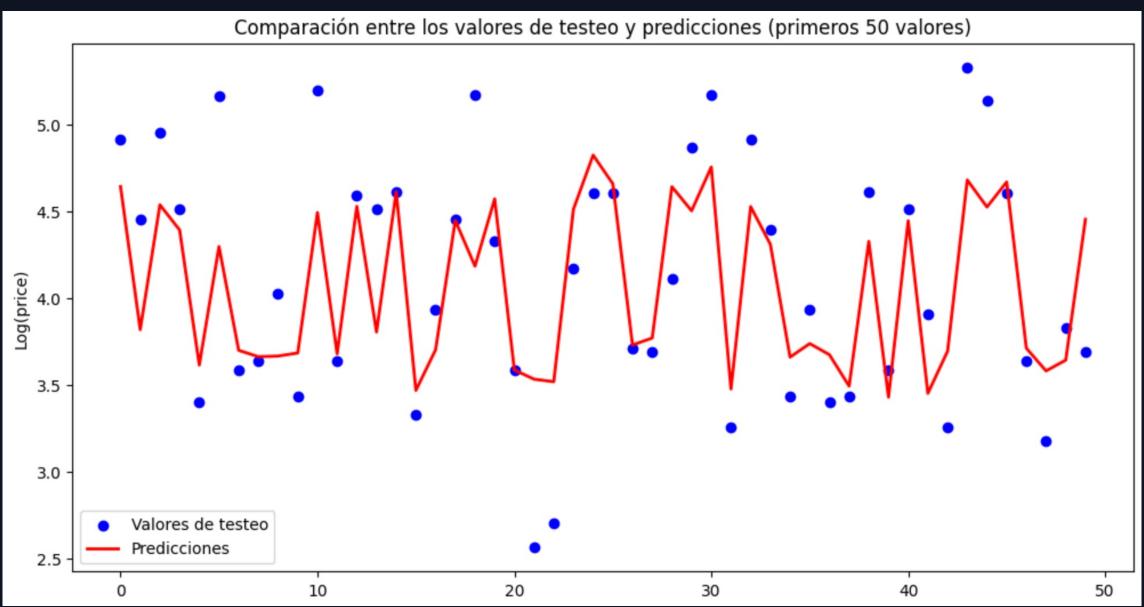
Con el dataset final se obtuvieron 2 conjuntos de datos adicionales: transformando todas las variables a su raíz cuadrada y transformando todas sus variables a su logaritmo. Esto con el final de reducir la escala de las variables haciendo así que los datos sean más uniformes.

Con estos 3 datasets se realizaron 2 modelos de regresión distintos: Regresión Lineal para el dataset transformado tanto en raíz cuadrada como en logaritmo y Regresión de Lasso Lars para el dataset con las variables sin transformar

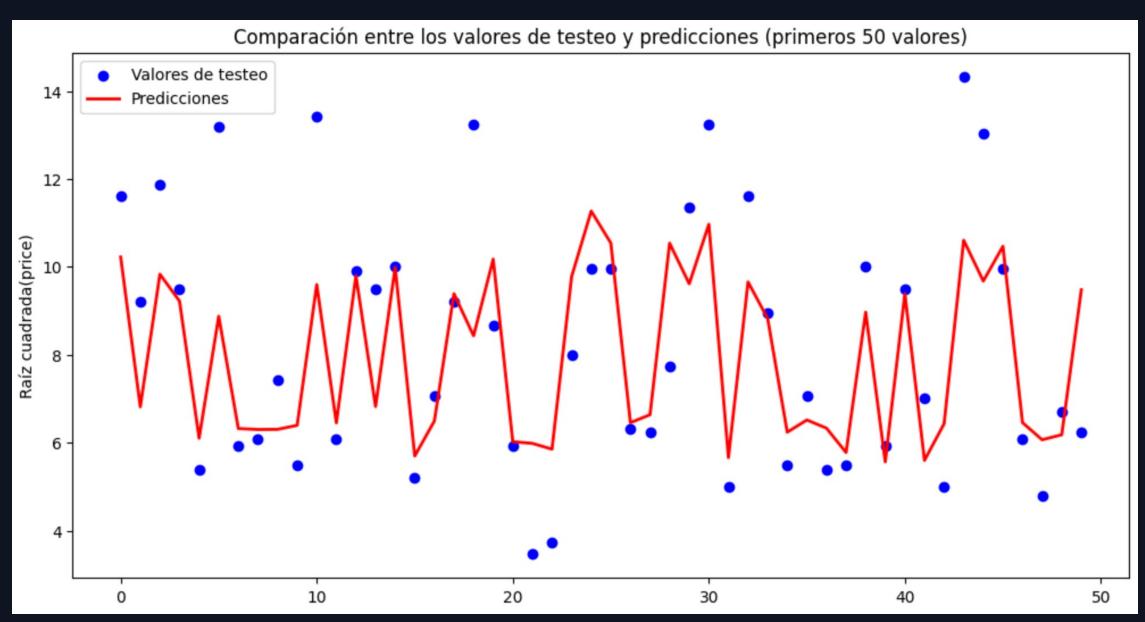
Modelo de Lasso Lars



Modelo de Regresión Lineal (log)

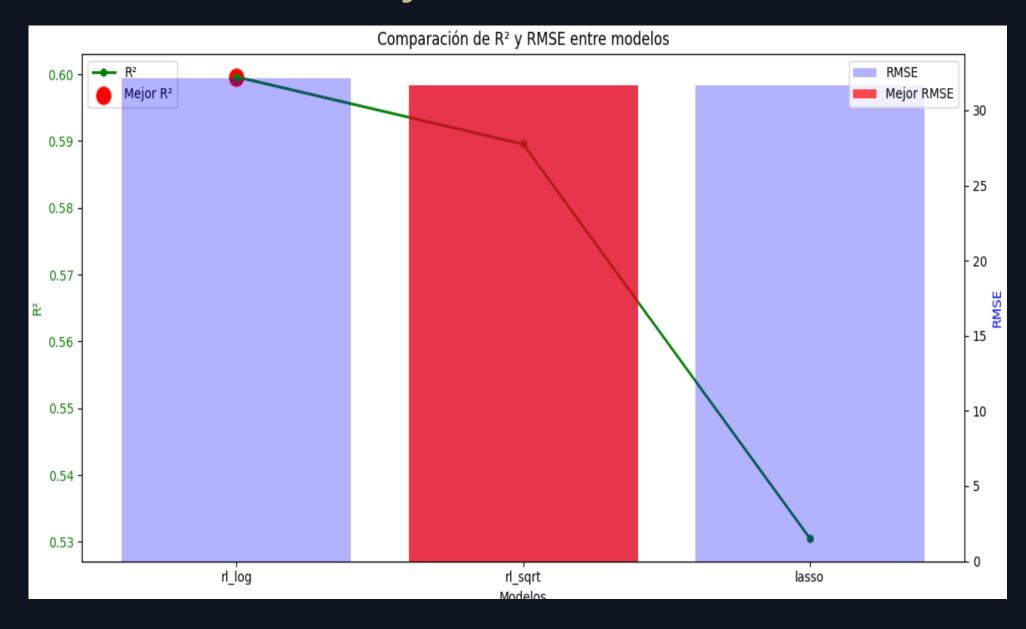


Modelo de Regresión Lineal (sqrt)



Elección de los mejores 2 modelos

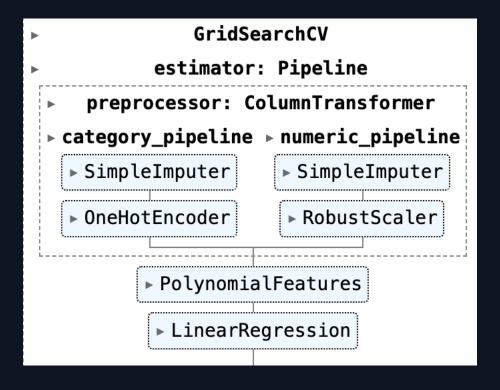
Se realizó una comparación entre los 3 modelos viendo su RMSE y su R². De esta manera se concluyó que los mejores 2 modelos fueron ambas regresiones lineales, una por obtener el mejor R² (0.6) y la otra por obtener el menor RMSE (31.7)



Mejora de los modelos

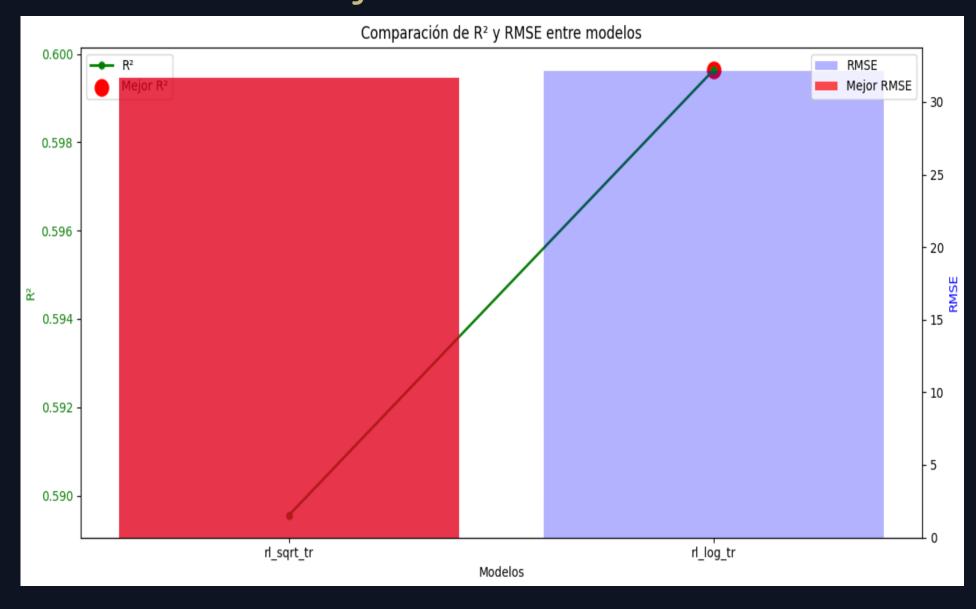
Dado que los mejores 2 modelos correspondían a una Regresión lineal, se realizó un preprocesador de ColumnTransformer para que los modelos sean capaces de encodear e imputar.

Además, se utilizó GridSearch en ambos modelos con el hiperparametro del grado del polinomio entre 1 y 4.



Elección de los mejores 2 modelos

Finalmente, con estas mejoras de ambos modelos, se obtuvieron métricas similares a las anteriores, por lo que el modelo escogido como final fue el que corresponde a las variables transformadas a su raíz cuadrada, ya que si bien no tiene el mejor R^2 (0.59 vs 0.60), tiene el RMSE más bajo (31.7 vs 32.2 euros), y lo que se busca en el modelo, además de un buen valor de R², es minimizar el error.



Conclusión

Luego de realizar varios modelos de regresión con 3 tipos de variables (originales, transformadas a raíz cuadrada y a logaritmo), se observó que el modelo con el mejor desempeño fue la Regresión Lineal con la raíz cuadrada de cada variable. Para escogerlo debió compararse con la otra Regresión Lineal, pero puesto que sus valores de R² eran similares, se optó por tomar la decisión según el RMSE, los cuales estaban en la misma unidad de medida (Euros), teniendo una diferencia de 0.5 euros.

Tomando en cuenta el gráfico de la comparación de los valores de testeo y las predicciones de este modelo, vemos que las predicciones se ajustan más a los valores de prueba cuando estos son intermedios, es decir, alojamientos que no poseen valores de rentas extremadamente altos ni bajos. Aunque el rendimiento del modelo es aceptable, existe un margen para mejoras, por lo que podría no solo servir como un complemento para tomar mejores decisiones a la hora de poner en renta o querer rentar un alojamiento, si no que podría ser un modelo más robusto, capaz de tomar decisiones optimas según los parámetros requeridos.