

## ## Introducción

Se comienza con dos dataset:

- Melbourne df

Obtenido de:

[https://cs.famaf.unc.edu.ar/~mteruel/datasets/diplodatos/melb\\_data.csv](https://cs.famaf.unc.edu.ar/~mteruel/datasets/diplodatos/melb_data.csv)  
' )

Consta de 13580 entradas y de las siguientes columnas:

'Suburb', 'Address', 'Rooms', 'Type', 'Price', 'Method', 'SellerG',  
'Date', 'Distance', 'Postcode', 'Bedroom2', 'Bathroom', 'Car',  
'Landsize', 'BuildingArea', 'YearBuilt', 'CouncilArea', 'Latitude',  
'Longitude', 'Regionname', 'Propertycount'.

- Airbnb df:

Obtenido de:

[https://cs.famaf.unc.edu.ar/~mteruel/datasets/diplodatos/cleansed\\_listings\\_decl8.csv](https://cs.famaf.unc.edu.ar/~mteruel/datasets/diplodatos/cleansed_listings_decl8.csv)

Consta de 22895 entradas y de las siguientes columnas:

'id', 'listing\_url', 'scrape\_id', 'last\_scraped', 'name', 'summary',  
'space', 'description', 'neighborhood\_overview', 'notes', 'transit',  
'access', 'interaction', 'house\_rules', 'picture\_url', 'host\_id',  
'host\_url', 'host\_name', 'host\_since', 'host\_location', 'host\_about',  
'host\_response\_time', 'host\_response\_rate', 'host\_is\_superhost',  
'host\_thumbnail\_url', 'host\_picture\_url', 'host\_neighborhood',  
'host\_verifications', 'host\_has\_profile\_pic', 'host\_identity\_verified',  
'street', 'neighborhood', 'city', 'suburb', 'state', 'zipcode',  
'smart\_location', 'country\_code', 'country', 'latitude', 'longitude',  
'is\_location\_exact', 'property\_type', 'room\_type', 'accommodates',  
'bathrooms', 'bedrooms', 'beds', 'bed\_type', 'amenities', 'price',  
'weekly\_price', 'monthly\_price', 'security\_deposit', 'cleaning\_fee',  
'guests\_included', 'extra\_people', 'minimum\_nights', 'maximum\_nights',  
'calendar\_updated', 'has\_availability', 'availability\_30',  
'availability\_60', 'availability\_90', 'availability\_365',  
'calendar\_last\_scraped', 'number\_of\_reviews', 'first\_review',  
'last\_review', 'review\_scores\_rating', 'review\_scores\_accuracy',  
'review\_scores\_cleanliness', 'review\_scores\_checkin',  
'review\_scores\_communication', 'review\_scores\_location',  
'review\_scores\_value', 'requires\_license', 'license',  
'instant\_bookable', 'cancellation\_policy',  
'require\_guest\_profile\_picture', 'require\_guest\_phone\_verification',  
'calculated\_host\_listings\_count', 'reviews\_per\_month'.

## ## Criterios de exclusión de ejemplos

Para excluir ejemplos dentro de las características numéricas utilizamos el método de los 3 sigmas: filtramos los valores por encima o por debajo de 3 desvíos estándar de la media de cada columna numérica. Vemos así que el valor de desviación estándar relativa (el desvío dividido la media) disminuye luego de aplicado este método. De esta manera eliminamos un 7.2% de los registros. Para el caso de las características categóricas lo que utilizamos es el criterio de filtrar aquellas categorías cuya participación es menor a un 0.5% del total de registros. De esta manera eliminamos un 5.6% de los registros.

## Características seleccionadas

### Características numéricas

Las características numéricas seleccionadas de Melbourne df:

- Rooms: Cantidad de habitaciones
- Bedroom2: Cantidad de dormitorios
- Bathrooms: Cantidad de baños
- Distance: Distancia al centro de la ciudad
- Landsize: Superficie del terreno
- BuildingSize: Superficie construida
- Car: Cantidad de cocheras
- YearBuilt: Año de construcción
- Postcode: Código Postal

El Precio para evaluar las predicciones.

Algunas características numéricas fueron agregadas de Airbnb df. Estas fueron vinculadas al df de Melbourne a través del Zipcode/Postcode. Se tomaron aquellos Zipcodes de Airbnb con más de 25 entradas. Sólo un 12% de Postcodes no fueron vinculados. Para vincular a un Zipcode varios registros las características fueron resumidas en su valor máximo, mínimo y promedio. Luego, las características numéricas agregadas de Airbnb df son:

- Price\_min: Precio mínimo de alquiler por día
- Price\_max: Precio máximo de alquiler por día
- Price\_mean: Precio promedio de alquiler por día
- Weekly\_price\_min: Precio mínimo de alquiler por semana
- Weekly\_price\_max: Precio máximo de alquiler por semana
- Weekly\_price\_mean: Precio de alquiler por semana
- Monthly\_price\_min: Precio mínimo de alquiler por mes
- Monthly\_price\_max: Precio máximo de alquiler por mes

- Monthly\_price\_mean: Precio promedio de alquiler por mes
- Zipcode: Código postal

### ### Características categóricas

Las características categóricas seleccionadas de Melbourne df:

- Type: tipo de propiedad. 3 valores posibles
- Suburb: Suburbio. 314 valores posibles
- CouncilArea: Municipio. 33 valores posibles

Y las características categóricas agregadas de Airbnb df:

- Neighborhood: Barrio. 59 valores posibles

### ### Transformaciones:

- Todas las características categóricas fueron codificadas con el método getdummies.
- Todas las características numéricas, menos el precio (la cual es la variable objetivo), fueron estandarizadas utilizando el método MinMaxScaler.
- Todas las características numéricas con valores nulos (\*) fueron imputadas utilizando el algoritmo KNN, con n\_neighbours = 5.

\*Cantidad de valores nulos por columna:

Car	111
BuildingArea	8783
YearBuilt	7179
CouncilArea	1485
zipcode	2270
airbnb_price_min	2270
airbnb_price_mean	2270
airbnb_price_max	2270
airbnb_weekly_price_min	10716
airbnb_weekly_price_mean	10716
airbnb_weekly_price_max	10716
airbnb_monthly_price_min	10899
airbnb_monthly_price_mean	10899
airbnb_monthly_price_max	10899

### ### Datos aumentados

- Se agregaron las 20 primeras columnas obtenidas a través del método de PCA, aplicado sobre el conjunto de datos totalmente procesado, menos la columna Price. Estas 20 columnas representan a un 65% de la varianza total.
- Se agregó al final del procesamiento la columna Price.

### ### Dataset final

El dataset final queda con 247 columnas numéricas y 18498 filas.