

En este documento se realiza el reporte de las operaciones que realizaron para obtener el conjunto de datos final, incluyendo:

1. Criterios de exclusión (o inclusión) de filas
2. Interpretación de las columnas presentes
3. Todas las transformaciones realizadas

Este documento es de uso técnico exclusivamente, y su objetivo es permitir que otros desarrolladores puedan reproducir los mismos pasos y obtener el mismo resultado.

## Criterios de exclusión de filas y columnas

Columnas del dataset que a priori no vamos a tener en cuenta para los análisis ya que no intervendrían en el valor de la propiedad

Address: se indica una dirección para cada propiedad, hay otras columnas que aportan información sobre la ubicación de la propiedad que permitirían un mejor agrupamiento por precio, por ejemplo, Regionname y CouncilArea.

SellerG: el nombre del vendedor no debería interferir con el valor de la propiedad

Date: la fecha de venta no sería indicativo del valor de la propiedad

Latitude y Longitude: podría pensarse a la longitud y latitud de la misma manera que a dirección (Address). Cabe aclarar que estas dos variables combinadas podrían ser útiles en cuanto a la ubicación específica de las propiedades en relación a otras características importantes de la región/barrio donde se encuentra (por ejemplo, cercanía a centros educativos, comercios etc).

Propertycount: es una variable relacionada a Suburb, se encuentra repetida con cada nombre de barrio diferente

Method: El método de venta no parece influir sobre el precio de las propiedades. El método SA parece tener valores más bajos pero representa un porcentaje de datos bajo en comparación al total de transacciones. según lo observado con anterioridad. Por otro lado el tipo de propiedad (Type) si parece influir en el precio.

Importante: La variable Postcode se puede utilizar como punto de unión entre datasets. Será utilizada posteriormente para incluir datos desde el dataset Airbnb.

A partir de los análisis descriptivos realizados, los criterios para la eliminación de filas del dataset fueron los siguientes:

- Se descartaron aquellas filas para las cuales Bedroom2 tenía un valor 0 (16) y valores mayores a 7 (20).

- Se eliminaron las filas para las cuales Bathroom tenía un valor 0 (34) y valores mayores a 5 (9).

- Se eliminaron las filas para las cuales Car tenía un valor mayor a 6 (21).

- Se eliminaron las filas para las cuales Buildingarea era igual a 0 (no consideramos propiedades vendidas al "pozo").

- Se eliminaron las filas para las cuales YearBuilt era menor al año 1850.

- Se eliminaron las filas para las cuales Councilarea era igual a "Macedon Ranges" y "Moorabool" (no pertenecen a Melbourne) y Unavailable.

## **Selección de variables:**

Variables categóricas:

- Suburb: Barrio donde se encuentra la propiedad

- Type: tipo de propiedad, 3 valores posibles,

- CouncilArea: nombre de la ciudad, 30 valores posibles.

- Regionname: Región geográfica donde se encuentra la propiedad, 6 valores posibles.

*Todas las características categóricas fueron codificadas con un método OneHotEncoding*

Variables numéricas:

- Rooms: Número de ambientes de la propiedad

- Price: Precio de la propiedad

- Distance: Distancia de la propiedad al centro

- Postcode: Código postal

- Bedroom2: Número de habitaciones recolectadas de un dataset secundario

- 64 .- Bathroom: Número de baños
- 65 .- Car: Número de cocheras
- 66 .- Landsize: Metros cuadrados de parque
- 67 .- BuildingArea: Metros cuadrados construidos
- 68 .- YearBuilt: Año de construcción.

69

70 Se seleccionaron las columnas "Price", "city" y "zipcode" del dataset Airbnb. City es  
71 equivalente a la columna CouncilArea y es utilizada para imputar a la misma. A su vez,  
72 Zipcode es el equivalente a Postcode y estas columnas son utilizadas para "unir" ambos  
73 datasets. El precio promedio de alquiler diario (Price) es la variable que se agrega al  
74 dataset, si bien es el que menos información contiene es el que menos valores faltantes  
75 presenta.

76

## 77 **Transformaciones realizadas:**

78

79 .- Se agrupan los valores correspondientes a Victoria de la variable Regioname para que  
80 esta región geográfica sea representativa en el dataset.

81 .- En el dataset Airbnb se reemplazaron los códigos postales con valores erróneos o mal  
82 anotados

83 .- Los valores faltantes de la columna CouncilArea fueron imputados a partir de los  
84 presentes en la variable City de Airbnb, utilizando Postocode y Zipcode como variables  
85 "clave". Los valores faltantes restantes fueron imputados utilizando la columna Suburb.

86 .- A partir del El precio promedio de alquiler diario (Price), se calculó el precio  
87 promedio de alquiler por ciudad, columna que fue agregada al dataset.

88 .- Las variables YearBuilt y BuildingArea fueron imputadas utilizando IterativeImputer  
89 con estimador KNeighborsRegressor. Para aplicar este método es necesario escalar los  
90 datos, dado que KNN utiliza como distancia entre valores la distancia euclidia, la cual  
91 necesita que los valores se encuentren en rangos equivalentes (se utiliza el método  
92 (PowerTransformer)).

93

## 94 **Componentes Principales (datos aumentados)**

95

96 Se genero un dataframe (df\_melb\_final) con el resultado de la codificación de las  
97 variables categóricas y el escalado e imputación de las variables numéricas. Luego, se  
98 procedió a realizar un análisis de componentes principales utilizando todas las  
99 variables disponibles (489).

100 Como resultado, se obtuvo que con los primeros 21 PCAs recuperamos  
101 aproximadamente el 90% de la variancia contenida en los datos. Estos PCAs fueron  
102 incluidos en un nuevo dataframe df\_melb\_pca\_final.to\_csv

103