

Andrés Felipe Rivas Luna

Daniel Felipe Galeano Tabares

Jorge Andrés Salas Barrera

LABORATORIO 3

El objetivo de este informe es proporcionar una guía detallada para entender y preparar el conjunto de datos "housing_fincaraiz.csv", siguiendo los lineamientos de IBM en lo que respecta a la comprensión de datos. Primero, se recolectaron los datos iniciales, identificando las variables clave como número de habitaciones, numero de baños, el precio, área construida. Luego, se describieron los datos, confirmando que están en formato .csv, con 31 columnas y cada variable con 8428 registros de datos, incluyendo datos numéricos, alfanuméricos, símbolos y texto. Finalmente, se verificó la calidad de los datos corrigiendo errores de escritura, ajustando nombres de variables y unidades de medida, y eliminando caracteres no deseados. Este proceso asegura que los datos estén limpios y listos para un análisis más profundo, proporcionando una base sólida para la toma de decisiones basada en datos.

Verificación de datos

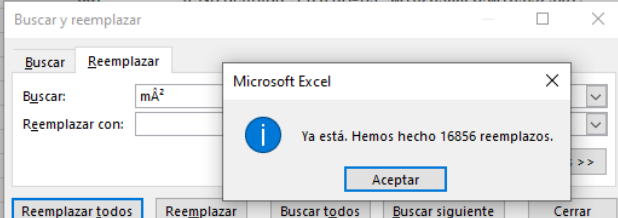
- Se corrigen los errores tipológicos de los nombres de cada variable

Porter	a / R Zonas Verde	Sal	Comu	Balc	Barra estilo	Calentador	Chimenea	Cit	fono
1	0	0	1	1	0	0	0	1	
1	0	0	1	1	1	1	0	0	
0	0	0	1	1	1	1	0	0	

- Para las columnas de área construida y área privada, vamos a generalizar los valores numéricos expuestos en m2, para ello vamos a eliminar “mÂ²” y vamos a especificar en el nombre de la variable que estos valores se encuentran en m2

area_construida	area_privada
92 mÂ²	92 mÂ²
56 mÂ²	56 mÂ²
144 mÂ²	144 mÂ²
31 mÂ²	31 mÂ²
52 mÂ²	52 mÂ²
150 mÂ²	150 mÂ²
110 mÂ²	100 mÂ²
53 mÂ²	47 mÂ²
111 mÂ²	0 mÂ²
264 mÂ²	264 mÂ²
97 mÂ²	0 mÂ²
87 mÂ²	82 mÂ²

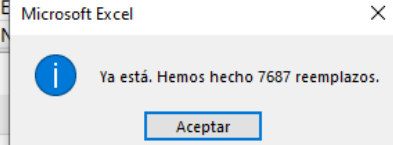
area_construida	area_privada	estrato	estado	antiguedad	administraci	precio_m2	Ascer
92	92	4	No definida	9 a 15 aA±os	\$A 622.000 C	\$A 6.521.739,	
56	56	6	No definida	1 a 8 aA±os	\$A 523.000 C	\$A 8.392.857	
144							
31							
52							
150							
110							
53							
111							
264							
97							
87	82	4	No definida	9 a 15 aA±os	\$A 350.000 C	\$A 4.712.643,	
175	175	4	No definida	16 a 30 aA±o	No definida	\$A 4.285.714,	
391	0	3	No definida	16 a 30 aA±o	No definida	\$A 1.994.884,	
218	0	4	Bueno	mÂ;s de 30 a	No definida	\$A 3.027.522,	
50	0	5	Bueno	9 a 15 aA±os	\$A 310.000 C	\$A 6.800.000	
140	140	4	Bueno	mÂ;s de 30 a	\$A 460.000 C	\$A 5.071.428,	
90	0	4	Bueno	mÂ;s de 30 a	\$A 369.000 C	\$A 4.000.000	
111	111	6	Bueno	9 a 15 aA±os	\$A 786.000 C	\$A 11.261.26,	



- Para la variable antigüedad corregiremos los errores tipográficos para ciertos datos, por ejemplo “años” por “años”

antigüedad
9 a 15 años
1 a 8 años
16 a 30 años
menor a 1 año
1 a 8 años
mas de 30 años
16 a 30 años
9 a 15 años

area_privada(estrato)	estado	antigüedad	administracion	precio_m2	Ascen
92	4 No definida	9 a 15 años	\$ 622.000 COP	\$ 6.521.739,13	
56	6 No definida	1 a 8 años	\$ 523.000 COP	\$ 8.392.857,14	
144	6 No definida	16 a 30 años	\$ 620.000 COP	\$ 6.597.222,22	
31	4 No definida				354,84
52	4 No definida				923,08



Buscar y reemplazar

Buscar: años

Reemplazar con: años

Opciones >>

antigüedad
9 a 15 años
1 a 8 años
16 a 30 años
menor a 1 año
1 a 8 años
mas de 30 años
16 a 30 años
9 a 15 años
16 a 30 años
mas de 30 años



Datos antigüedad formateados

- En las variables administración y precio m2 se va a eliminar los caracteres “\$” y todos aquellos caracteres que pueden afectar en la lectura de datos

administracion	precio_m2
\$ 622.000 COP	\$ 6.521.739,13 *m²
\$ 523.000 COP	\$ 8.392.857,14 *m²
\$ 620.000 COP	\$ 6.597.222,22 *m²
\$ 130.000 COP	\$ 7.419.354,84 *m²
\$ 219.000 COP	\$ 5.576.923,08 *m²
\$ 872.000 COP	\$ 6.533.333,33 *m²



administracion(COP)	precio_m2
622.000	6.521.739,13
523.000	8.392.857,14
620.000	6.597.222,22
130.000	7.419.354,84
219.000	5.576.923,08
872.000	6.533.333,33
135.000	4.181.818,18
125.000	3.207.547,17
No definida	3.873.873,87

- Por último, se verifican los datos de la variable ubicación y se corrigen los errores tipográficos

ubicacion	▼
Centro Internacional	
Calleja Baja	
Cerros de Suba	
Mazuren	
El plan	
La Cabrera	
Ciudad jardin sur	
Ub. industrial las delicias	
Pradera Norte	
Nicolas de federman	
Pontevedra	
Los cedros oriental	
Cedro Golf	
Sosiego sur	
Normandía	
Contador	
Quinta Paredes	
Cedro Golf	
La Cabrera	
Los Cerezos	

CONCLUSIONES

La verificación de la calidad de los datos es esencial para asegurar que el análisis posterior sea preciso y significativo. Al seguir estos pasos, se garantiza que los datos estén limpios y listos para un análisis más profundo, proporcionando así una base sólida para la toma de decisiones basada en datos.