

## Integrantes :

Paula Rios - 91105

Oscar Velasquez - 91273

## Laboratorio 3

### ● Recolectar datos iniciales

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
habitaciones	baños	parqueadero	area_construida	area_privada	estrato	estado	antigüedad	administración	precio_m2	Ascensor	Circuito cerrado	Parqueadero	Portería	Zonas Verdes	Salón Común	Balcón
2	2	1	92 m <sup>2</sup>	92 m <sup>2</sup>		4 No definida	9 a 15 años	\$Á 622.000 C	\$Á 6.521.739	1	1	1	1	0	0	
1	2	1	56 m <sup>2</sup>	56 m <sup>2</sup>		6 No definida	1 a 8 años	\$Á 523.000 C	\$Á 8.392.857	1	0	1	1	0	0	
3	4	2	144 m <sup>2</sup>	144 m <sup>2</sup>		6 No definida	16 a 30 años	\$Á 620.000 C	\$Á 6.597.222	1	0	0	0	0	0	
1	1	0	31 m <sup>2</sup>	31 m <sup>2</sup>		4 Excelente	menor a 1 año	\$Á 130.000 C	\$Á 7.419.354	1	1	1	0	0	0	
3	2	1	52 m <sup>2</sup>	52 m <sup>2</sup>		4 No definida	1 a 8 años	\$Á 219.000 C	\$Á 5.576.923	1	1	1	0	0	0	
3	3	1	150 m <sup>2</sup>	150 m <sup>2</sup>		6 Bueno	más de 30 años	\$Á 872.000 C	\$Á 6.533.333	1	1	1	0	0	0	
3	2	1	110 m <sup>2</sup>	100 m <sup>2</sup>		3 No definida	16 a 30 años	\$Á 135.000 C	\$Á 4.181.818	0	1	0	0	1	0	
3	2	0	53 m <sup>2</sup>	47 m <sup>2</sup>		3 No definida	9 a 15 años	\$Á 125.000 C	\$Á 3.207.547	0	0	1	1	1	1	
3	3	1	111 m <sup>2</sup>	0 m <sup>2</sup>		4 Excelente	16 a 30 años	No definida	\$Á 3.873.873	0	1	1	0	0	1	
4	4	2	264 m <sup>2</sup>	264 m <sup>2</sup>		5 Bueno	más de 30 años	\$Á 836.000 C	\$Á 5.303.030	1	0	0	0	0	0	
3	2	2	97 m <sup>2</sup>	0 m <sup>2</sup>		4 No definida	9 a 15 años	\$Á 272.000 C	\$Á 6.175.257	1	1	0	0	0	0	
2	3	2	87 m <sup>2</sup>	82 m <sup>2</sup>		4 No definida	9 a 15 años	\$Á 350.000 C	\$Á 4.712.643	1	1	1	1	0	1	
3	4	2	175 m <sup>2</sup>	175 m <sup>2</sup>		4 No definida	16 a 30 años	No definida	\$Á 4.285.714	0	0	0	0	1	0	
10	5	2	391 m <sup>2</sup>	0 m <sup>2</sup>		3 No definida	16 a 30 años	No definida	\$Á 1.994.884	0	0	0	0	0	0	
6	4	3	218 m <sup>2</sup>	0 m <sup>2</sup>		4 Bueno	más de 30 años	No definida	\$Á 3.027.522	0	1	0	0	0	0	
1	2	1	50 m <sup>2</sup>	0 m <sup>2</sup>		5 Bueno	9 a 15 años	\$Á 310.000 C	\$Á 6.800.000	0	0	0	1	0	0	
3	3	1	140 m <sup>2</sup>	140 m <sup>2</sup>		4 Bueno	más de 30 años	\$Á 460.000 C	\$Á 5.071.428	0	1	0	1	0	1	
4	3	1	90 m <sup>2</sup>	0 m <sup>2</sup>		4 Bueno	más de 30 años	\$Á 369.000 C	\$Á 4.000.000	1	0	1	1	0	1	
2	3	2	111 m <sup>2</sup>	111 m <sup>2</sup>		6 Bueno	9 a 15 años	\$Á 786.000 C	\$Á 11.261.26	1	1	1	0	0	0	
12	5	0	375 m <sup>2</sup>	320 m <sup>2</sup>		3 Bueno	16 a 30 años	No definida	\$Á 1.493.333	0	0	0	0	0	0	
2	3	2	111 m <sup>2</sup>	73 m <sup>2</sup>		4 No definida	16 a 30 años	\$Á 483.000 C	\$Á 3.423.423	1	0	1	1	0	1	

- Las variables más prometedoras son: precio, habitación, baños, nombre, area\_construida y area\_privada
- Fuente única de datos, no se está fusionando con otra fuente
- No hay datos faltantes, los datos están completos

### ● Describir los datos

- El archivo es en formato .csv
- Los datos se obtuvieron mediante web scraping
- El archivo tiene 31 columnas (atributos) y 8428 filas (registros)
- Contiene datos numéricos, alfanuméricos y de texto

### ● Explorar los datos

Las variables que son prometedoras para un análisis más profundo son: antigüedad, area\_construida, area\_privada, nombre y ubicación.

### ● Verificar los datos

Errores tipográficos en los datos:

- Corrección nombre de atributos

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
1	habitacion	baños	parquead	area_construida	area_privada	estrato	estado	antigüedad	administración	precio_m2	Ascensor	Circuito cerrado	Parqueadero	Portería	Zonas Verdes	Salón Común	Balcón
2	2	2	1	92 m <sup>2</sup>	92 m <sup>2</sup>		4 No definida	9 a 15 años	\$Á 622.000 C	\$Á 6.521.739	1	1	1	1	0	0	

- Se cambian nombres de variables `area_construida` y `area_privada`, por `area_construida_m2` y `area_privada_m2`, se eliminan las unidades(m<sup>2</sup>) de los registros de las variables `precio_m2`, `area_construida_m2` y `area_privada_m2`

The screenshot shows an Excel spreadsheet with a table of housing data. The columns are labeled: habitacion, baños, parquead, area\_con, area\_priv, estrato, estado, and antigüedad. A 'Buscar y reemplazar' dialog box is open, showing 'Buscar: m²' and 'Reemplazar con:'. A message box says 'Ya está. Hemos hecho 25284 reemplazos.'

- Se cambia nombre de la variable `antigüedad` por `antigüedad_años`, se elimina la unidad de la antigüedad (años, años)

The screenshot shows an Excel spreadsheet with a table of housing data. The columns are labeled: A, B, C, D, E, F, G, H. A 'Buscar y reemplazar' dialog box is open, showing 'Buscar: años' and 'Reemplazar con:'. A message box says 'Ya está. Hemos hecho 298 reemplazos.'

- Se cambian los registros de la variable `antigüedad_años` que contengan la palabra "mas" por la palabra "mas"

Microsoft Excel interface showing a search and replace operation. The "Buscar y reemplazar" dialog box is open, with "Buscar" set to "mas de 30" and "Reemplazar con" set to "mas". The "Reemplazar todos" button is highlighted. A message box states: "Ya está. Hemos hecho 1643 reemplazos." The spreadsheet data includes columns for habitación, baños, parquead, area\_cons, area\_priv, estrato, estado, and antigüedad\_años.

habitación	baños	parquead	area_cons	area_priv	estrato	estado	antigüedad_años
2	2	1	92	92	4	No definida	9 a 15
1	2	1	56	56	6	No definida	1 a 8
3	4	2	144	144	6	No definida	16 a 30
1	1	0	31	31	4	Excelente	menor a 1
3	2	1	52	52	4	No definida	9 a 15
3	3	1	150	150	6	Bueno	mas de 30
3	2	1	110	100	3	No definida	9 a 15
3	2	0	53	47	3	No definida	16 a 30
3	3	1	111	0	4	Excelente	menor a 1
4	4	2	264	264	5	Bueno	mas de 30
3	2	2	97	0	4	No definida	9 a 15
2	3	2	87	82	4	No definida	9 a 15
3	4	2	175	175	4	No definida	16 a 30
10	5	2	391	0	3	No definida	16 a 30
6	4	3	218	0	4	Bueno	mas de 30
1	2	1	50	0	5	Bueno	9 a 15
3	3	1	140	140	4	Bueno	mas de 30
4	3	1	90	0	4	Bueno	mas de 30
2	3	2	111	111	6	Bueno	9 a 15
12	5	0	375	320	3	Bueno	16 a 30
2	3	2	111	73	4	No definida	16 a 30

- Se elimina "\$Â" de la variable administracion y precio\_m2

Microsoft Excel interface showing a search and replace operation. The "Buscar y reemplazar" dialog box is open, with "Buscar" set to "\$Â" and "Reemplazar con" set to an empty field. The "Reemplazar todos" button is highlighted. A message box states: "Ya está. Hemos hecho 15008 reemplazos." The spreadsheet data includes columns for habitación, baños, parquead, area\_cons, area\_priv, estrato, estado, antigüedad\_años, and precio\_m2.

habitación	baños	parquead	area_cons	area_priv	estrato	estado	antigüedad_años	precio_m2
2	2	1	92	92	4	No definida	9 a 15	622.000 COP
1	2	1	56	56	6	No definida	1 a 8	523.000 COP
3	4	2	144	144	6	No definida	16 a 30	620.000 COP
1	1	0	31	31	4	Excelente	menor a 1	130.000 COP
3	2	1	52	52	4	No definida	9 a 15	5.576.923,08
3	3	1	150	150	6	Bueno	mas de 30	6.533.333,33
3	2	1	110	100	3	No definida	9 a 15	4.181.818,18
3	2	0	53	47	3	No definida	16 a 30	3.207.547,17
3	3	1	111	0	4	Excelente	menor a 1	3.873.873,87
4	4	2	264	264	5	Bueno	mas de 30	5.303.030,3*
3	2	2	97	0	4	No definida	9 a 15	272.000 COP
2	3	2	87	82	4	No definida	9 a 15	350.000 COP
3	4	2	175	175	4	No definida	16 a 30	No definida
10	5	2	391	0	3	No definida	16 a 30	No definida
6	4	3	218	0	4	Bueno	mas de 30	No definida
1	2	1	50	0	5	Bueno	9 a 15	310.000 COP
3	3	1	140	140	4	Bueno	mas de 30	460.000 COP
4	3	1	90	0	4	Bueno	mas de 30	369.000 COP
2	3	2	111	111	6	Bueno	9 a 15	786.000 COP
12	5	0	375	320	3	Bueno	16 a 30	No definida
2	3	2	111	73	4	No definida	16 a 30	483.000 COP

- Se elimina el "\*" de la variable precio\_m2, insertando una nueva columna en la que estaran unicamente valores numericos, y se elimina la columna que contiene el "\*"

=SUSTITUIR(K2;"*";" ")								
	D	E	F	G	H	I	J	K
	area_cons	area_priv	estrato	estado	antigüedad_años	administracion	precio_m2	precio_m2
1	92	92	4	No definida	9 a 15	622.000 COP	6.521.739,13	6.521.739,13*
1	56	56	6	No definida	1 a 8	523.000 COP	8.392.857,14	8.392.857,14*
2	144	144	6	No definida	16 a 30	620.000 COP	6.597.222,22	6.597.222,22*
0	31	31	4	Excelente	menor a 1	130.000 COP	7.419.354,84	7.419.354,84*
1	52	52	4	No definida	1 a 8	219.000 COP	5.576.923,08	5.576.923,08*
1	150	150	6	Bueno	mas de 30	872.000 COP	6.533.333,33	6.533.333,33*
1	110	100	3	No definida	16 a 30	135.000 COP	4.181.818,18	4.181.818,18*
0	53	47	3	No definida	9 a 15	125.000 COP	3.207.547,17	3.207.547,17*
1	111	0	4	Excelente	16 a 30	No definida	3.873.873,87	3.873.873,87*
2	264	264	5	Bueno	mas de 30	836.000 COP	5.303.030,3	5.303.030,3*
2	97	0	4	No definida	9 a 15	272.000 COP	6.175.257,73	6.175.257,73*
2	87	82	4	No definida	9 a 15	350.000 COP	4.712.643,68	4.712.643,68*
2	175	175	4	No definida	16 a 30	No definida	4.285.714,29	4.285.714,29*
2	391	0	3	No definida	16 a 30	No definida	1.994.884,91	1.994.884,91*
3	218	0	4	Bueno	mas de 30	No definida	3.027.522,94	3.027.522,94*
1	50	0	5	Bueno	9 a 15	310.000 COP	6.800.000	6.800.000*
1	140	140	4	Bueno	mas de 30	460.000 COP	5.071.428,57	5.071.428,57*
1	90	0	4	Bueno	mas de 30	369.000 COP	4.000.000	4.000.000*
2	111	111	6	Bueno	9 a 15	786.000 COP	11.261.261,26	11.261.261,26*
0	375	320	3	Bueno	16 a 30	No definida	1.493.333,33	1.493.333,33*
2	111	73	4	No definida	16 a 30	483.000 COP	3.423.423,42	3.423.423,42*

6.521.739,13

	D	E	F	G	H	I	J	K	L
	area_cons	area_priv	estrato	estado	antigüedad_años	administracion	precio_m2	Ascensor	Circuito c
1	92	92	4	No definida	9 a 15	622.000 COP	6.521.739,13	1	1
1	56	56	6	No definida	1 a 8	523.000 COP	8.392.857,14	1	0
2	144	144	6	No definida	16 a 30	620.000 COP	6.597.222,22	1	0
0	31	31	4	Excelente	menor a 1	130.000 COP	7.419.354,84	1	1
1	52	52	4	No definida	1 a 8	219.000 COP	5.576.923,08	1	1
1	150	150	6	Bueno	mas de 30	872.000 COP	6.533.333,33	1	1
1	110	100	3	No definida	16 a 30	135.000 COP	4.181.818,18	0	1
0	53	47	3	No definida	9 a 15	125.000 COP	3.207.547,17	0	0
1	111	0	4	Excelente	16 a 30	No definida	3.873.873,87	0	1
2	264	264	5	Bueno	mas de 30	836.000 COP	5.303.030,3	1	0
2	97	0	4	No definida	9 a 15	272.000 COP	6.175.257,73	1	1
2	87	82	4	No definida	9 a 15	350.000 COP	4.712.643,68	1	1
2	175	175	4	No definida	16 a 30	No definida	4.285.714,29	0	0
2	391	0	3	No definida	16 a 30	No definida	1.994.884,91	0	0
3	218	0	4	Bueno	mas de 30	No definida	3.027.522,94	0	1
1	50	0	5	Bueno	9 a 15	310.000 COP	6.800.000	0	0
1	140	140	4	Bueno	mas de 30	460.000 COP	5.071.428,57	0	1
1	90	0	4	Bueno	mas de 30	369.000 COP	4.000.000	1	0
2	111	111	6	Bueno	9 a 15	786.000 COP	11.261.261,26	1	1
0	375	320	3	Bueno	16 a 30	No definida	1.493.333,33	0	0
2	111	73	4	No definida	16 a 30	483.000 COP	3.423.423,42	1	0

- Se eliminan errores tipograficos en la variable ubicacion

Chimeneá	Citofono	Cocina Int	Terraza	Vigilancia	Parques c	Estudio	Patio	Deposito	nombre	ubicacion	precio
0	1	0	1	0	0	0	0	0	0 Apartamento	Centro Internacional	600000000
0	0	0	1	0	0	0	0	0	0 Apartamento	Calleja Baja	470000000
0	0	0	0	0	0	0	0	0	0 Apartamento	Cerros de Suba	950000000
1	1	0	0	0	0	0	0	0	0 Apartamento	Mazuren	230000000
0	1	0	0	0	0	0	0	0	0 Apartamento	El plan	290000000
1	1	1	0	0	0	0	0	0	0 Apartamento	La Cabrera	980000000
0	0	1	0	0	1	1	1	1	0 Apartamento	Ciudad jardin sur	460000000
0	0	0	0	1	0	1	0	0	0 Apartamento	Ub. industrial las delicias	170000000
0	0	1	0	0	0	0	0	0	0 Casa	Pradera Norte	430000000
1	0	0	1	0	0	1	0	0	0 Apartamento	Nicolas de federman	1400000000
0	1	0	0	0	0	0	0	0	0 Apartamento	Pontevedra	599000000
1	1	0	0	1	0	0	0	0	0 Apartamento	Los cedros oriental	410000000
1	0	1	0	0	0	0	0	0	0 Casa	Cedro Golf	750000000
0	0	0	0	0	0	0	0	0	0 Casa	Sosiego sur	780000000
0	0	1	0	0	0	1	0	0	0 Casa	Normandia	660000000
0	0	1	0	0	0	0	0	0	1 Apartamento	Contador	340000000
1	1	1	1	1	0	1	0	0	0 Apartamento	Quinta Paredes	710000000
1	1	1	1	0	0	0	0	0	1 Apartamento	Cedro Golf	360000000
0	1	0	0	0	0	0	0	0	0 Apartamento	La Cabrera	1250000000
0	0	1	1	0	1	0	1	1	0 Casa	Los Cerezos	560000000
1	1	0	0	1	0	0	0	0	0 Apartamento	Cedro Golf	380000000

A continuación se anexa el archivo .csv luego de la limpieza de datos:

[housing\\_fincaraiz\\_Lim](#)