



DATA WOMEN'S COMPANY

W H E R E D A T A M A K E S S E N S E



Diana Maria Toro López



María de Lluch Gual Pérez



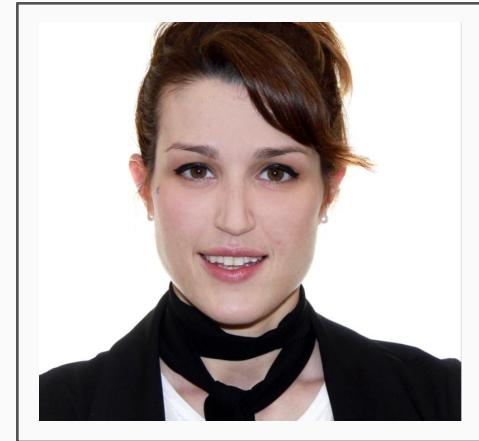
Sanja Aleksova



Sandra Julieth Castaño Reina



Elsa Cembrero Bonet



Mireia Hernández Lozano

ÍNDICE

1. INTRODUCCIÓN Y OBJETIVOS
2. DEFINIENDO EL DATA SET
3. NUESTRO SERVICIO DE CONSULTORÍA
4. ARQUITECTURA Y VALIDACIÓN DE LOS DATOS
5. VISUALIZACIÓN DE LAS MÉTRICAS
6. PRE-PROCESAMIENTO Y MODELADO
7. RESULTADOS
8. CONCLUSIONES



DATA WOMEN'S COMPANY

1. INTRODUCCIÓN Y OBJETIVOS



DATA WOMEN'S COMPANY

2. DEFINIENDO EL DATASET

opendatasoft EXPLORE MAP BUILDER API CHART BUILDER Login

14,780 records

Active filters [Clear all](#)

Text search Madrid

Filters

Madrid [🔍](#)

Host Response Time

within an hour	7,905
within a few hours	2,872
within a day	1,802
a few days or more	302

Host Response Rate

100	9,670
90	340
96	297
99	260
75	227
98	213
> More	

Host Verifications

Airbnb - Listings

Information Table Map Analyze Images Export API

This dataset is licensed under : CC 0 1.0

Flat file formats

- CSV [Whole dataset](#) [Only the 14780 selected records](#)
CSV uses semicolon (;) as a separator.
- JSON [Whole dataset](#) [Only the 14780 selected records](#)
- Excel [Whole dataset](#) [Only the 14780 selected records](#)

Geographic file formats

- GeoJSON [Whole dataset](#) [Only the 14780 selected records](#)
- Shapefile [Whole dataset](#) [Only the 14780 selected records](#)
⚠ This export format is limited to 50,000 records. You can download a smaller part of the dataset by filtering it.
- KML [Whole dataset](#) [Only the 14780 selected records](#)

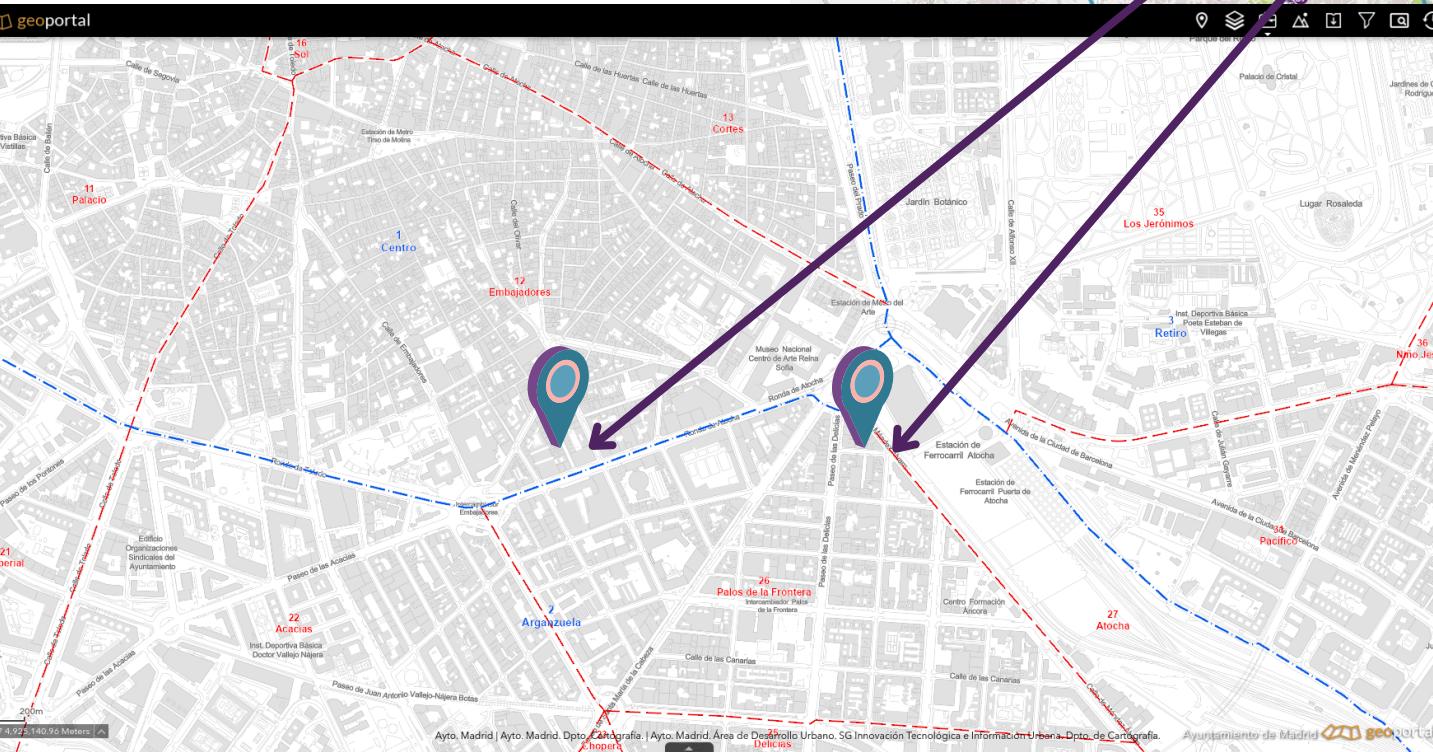
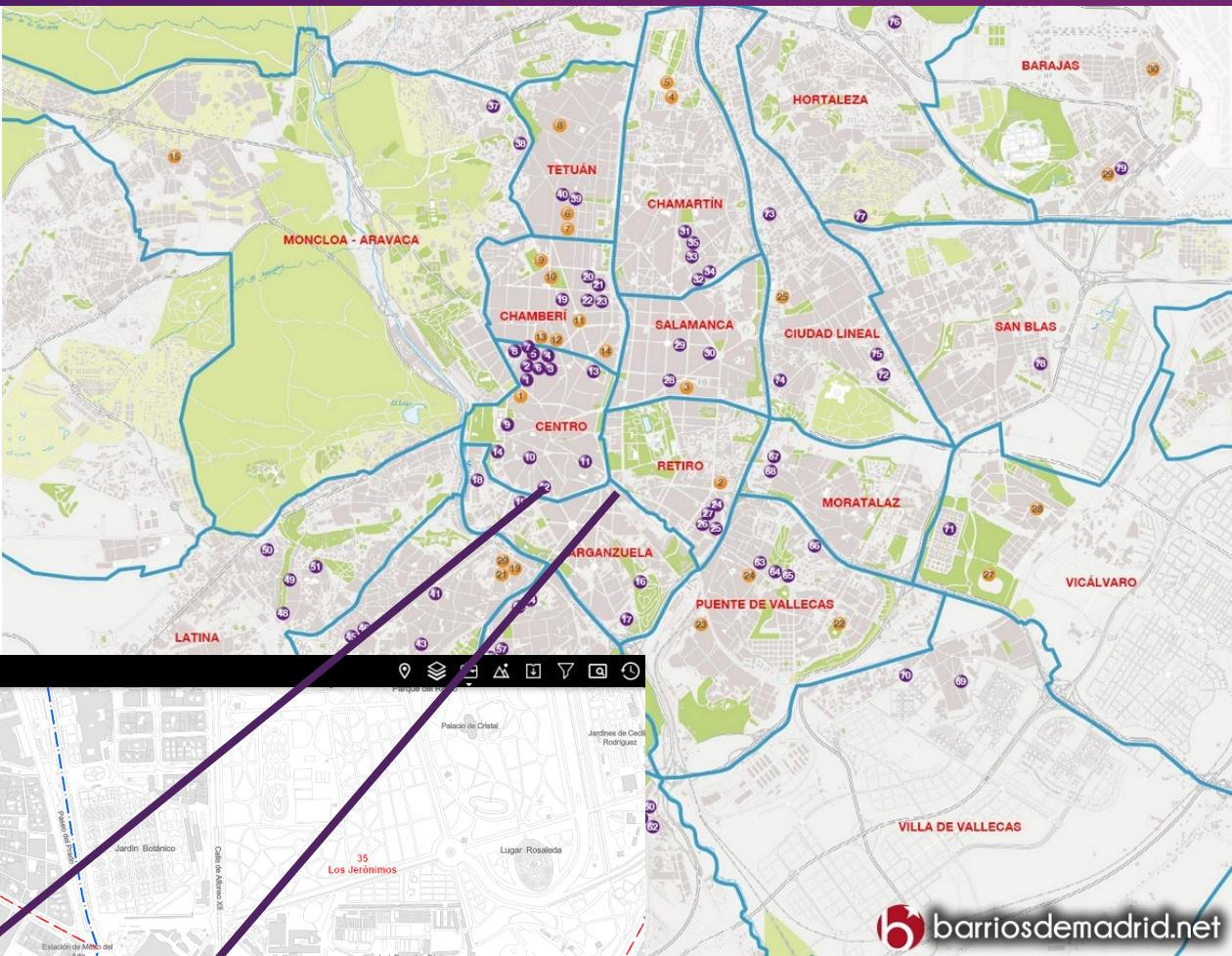
Data columns (total 89 columns):		
#	Column	Non-Null Count Dtype
0	ID	14780 non-null int64
1	Listing Url	14780 non-null object
2	Scrape ID	14780 non-null int64
3	Last Scrapped	14780 non-null object
4	Name	14779 non-null object
5	Summary	14189 non-null object
6	Space	10888 non-null object
7	Description	14774 non-null object
8	Experiences Offered	14780 non-null object
9	Neighborhood Overview	9134 non-null object
10	Notes	5644 non-null object
11	Transit	9066 non-null object
12	Access	8318 non-null object
13	Interaction	8228 non-null object
14	House Rules	9619 non-null object
15	Thumbnail Url	11960 non-null object
16	Medium Url	11960 non-null object
17	Picture Url	14751 non-null object
18	XL Picture Url	11960 non-null object
19	Host ID	14780 non-null int64
...		
87	Geolocation	14780 non-null object
88	Features	14779 non-null object

dtypes: float64(23), int64(13), object(53)
memory usage: 10.0+ MB

DATA W MEN'S COMPANY

3. NUESTRO SERVICIO DE CONSULTORÍA

- ¿A qué **precio base** debería publicar los apartamentos?
- ¿Cómo puede variar el **precio base** según la **cantidad de reviews**?
- Grado de implicación:
¿Es rentable ser **superanfitrión**?



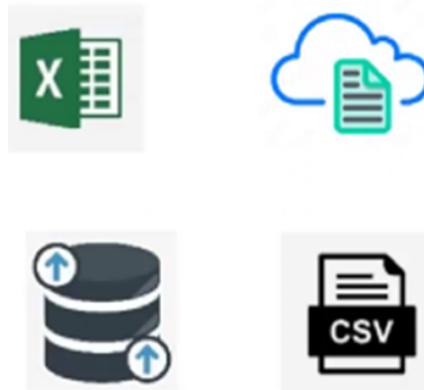
- Apartamentos
- Completos
- 2 habitaciones
- 2 camas
- 1 baño
- Hasta 6 huéspedes
- 1 incluido en precio

4. ARQUITECTURA Y VALIDACIÓN DE LOS DATOS



Procesos ETL: Extracción

FUENTE DE DATOS



- Analizar
- Interpretar
- Convertir
- **Cuidado y cautela**
- Horarios para realizar el proceso

Procesos ETL: Transformación

Transformación



- Aplicar regla de negocio
- Declarativas
- Independientes
- Inteligibles
- Útiles
- Practico al momento de transformar datos
- Seleccionar columnas
- Traducir códigos
- Codificar valores libres
- Valores calculados
- Unir datos de múltiples fuentes
- Generar campos claves
- Transponer o pivotar
- Dividir columnas en varias

DATOS DE DESTINO

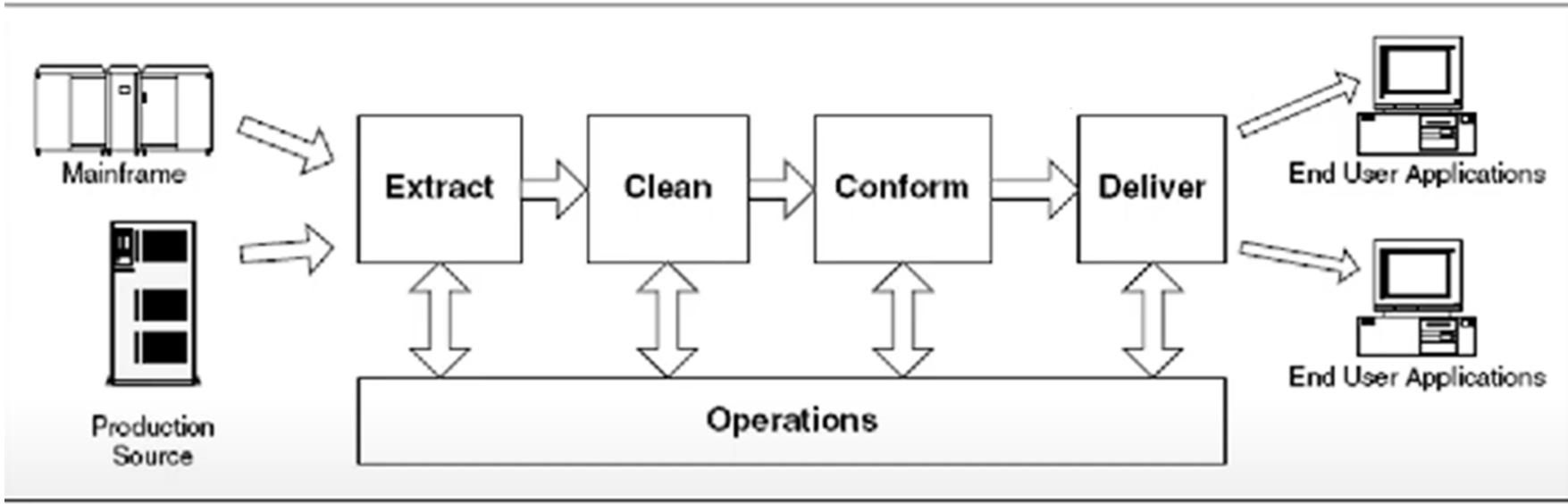
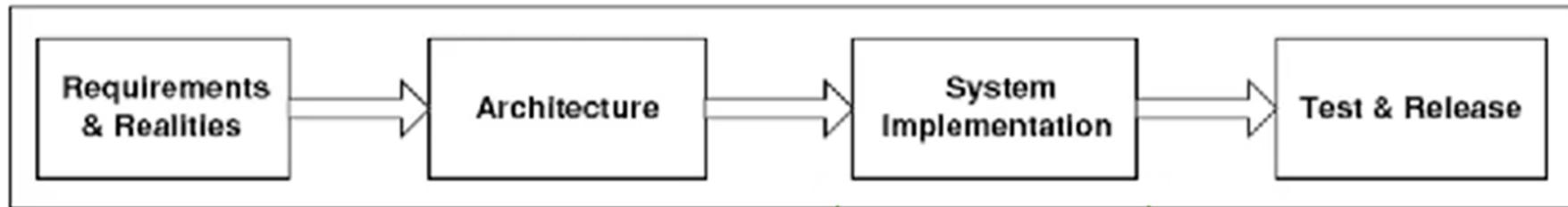


Procesos ETL: Carga

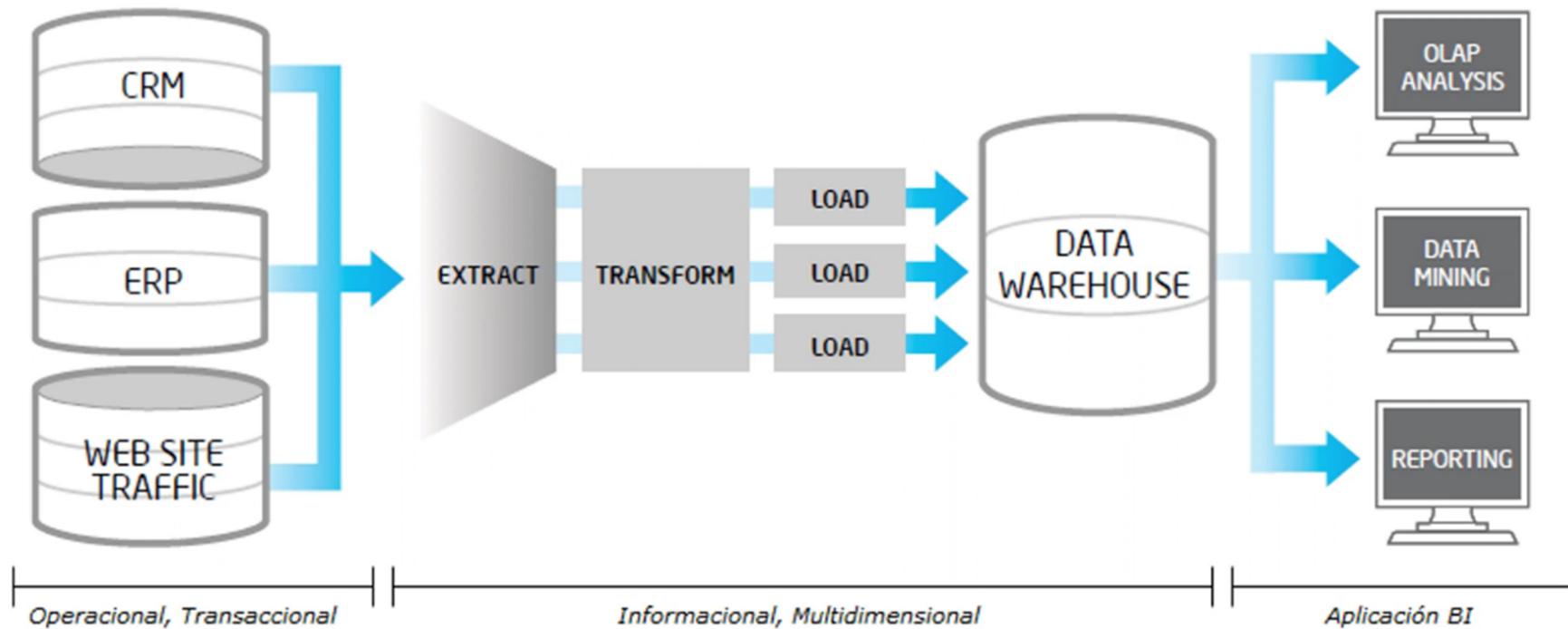
- Carga de datos al sistema de destino
- Políticas de la organización al cargar datos
- **Procesos de carga**
 - Acumulación simple
 - Rolling



DATA W MEN'S COMPANY



DATA W MEN'S COMPANY





5. VISUALIZACIÓN DE LAS MÉTRICAS

- Selección de 4 constantes:
 - Apartamento
 - Precio por día
 - 2 Habitaciones
 - Distritos
- Visualización de un mapa con los dos distritos y la mediana del precio por día
- 3 gráficos de barras con variables que afectan al precio
 - Afecta positivamente:
 - Ser superhost
 - El barrio
 - Afecta negativamente:
 - Número de reviews

6. PRE-PROCESAMIENTO Y MODELADO

- Selección -> características de los alojamientos

```
{r}
df_amarito <- dplyr:: filter(df, Neighbourhood.Cleansed == "embajadores", Property.Type == "apartment", Room.Type == "entire home/apt")
df_amarito
```

```
{r}
df_alvarito <- dplyr:: filter(df, Neighbourhood.Cleansed == "palos de moguer", Property.Type == "apartment", Room.Type == "entire home/apt")
df_alvarito
```



- Detección de outliers - características

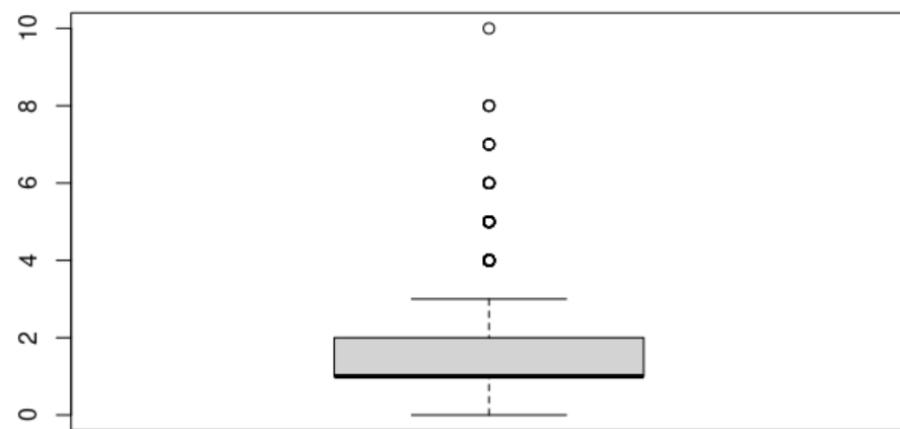
```
{r}
boxplot.stats(df_amaro$Bedrooms)
boxplot(df_amaro$Bedrooms)
```



- Nueva selección con rangos de outliers

```
{r}
amarito_s <- dplyr:: filter(df_amarito, Bedrooms <=3, Beds <= 5,
Number.of.Reviews <= 121, Price <=140, Accommodates <=7)
amarito_s
```

```
{r}
alvarito_s <- dplyr:: filter(df_alvarito, Bedrooms <=3, Beds <= 4,
Guests.Included <=3, Number.of.Reviews <= 49, Price <=130,
Accommodates <=8)
alvarito_s
```



DATA W[♂] MEN'S COMPANY

- Correlacionar -> relaciones entre variables
 - Gráfica y análisis de correlación
 - Generación de modelo con todas las variables observando las que no son útiles
 - Comprobación mediante la función **step ()** que selecciona las variables previstas
 - Obtención de fórmulas

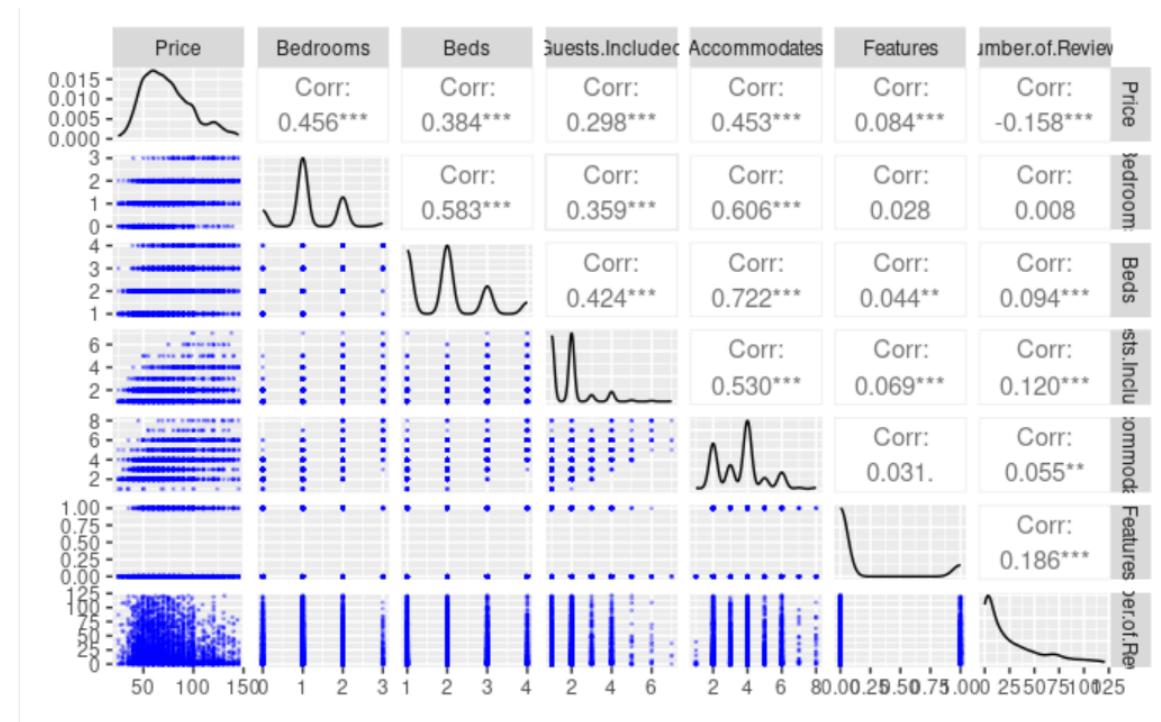
```
{r}
round(cor(x = amparito_s[,c("Price", "Bedrooms", "Beds", "Guests.Included", "Accommodates", "Features", "Number.of.Reviews")]), method = "pearson"), 3)
```

```
{r}
round(cor(x = amparito_s[,c("Price", "Bedrooms", "Beds", "Guests.Included", "Accommodates", "Features", "Number.of.Reviews")]), method = "spearman"), 3)
```

```
[r] ggpairs(amparo_s[,c("Price", "Bedrooms", "Beds", "Guests.Included", "Accommodates", "Features", "Number.of.Reviews")], lower = list(continuous = wrap("points", alpha = 0.3, size=0.3, color=blue)))
```

```
{r}
model_amarquito_s<- lm(data=amarquito_s, formula = Price~Bedrooms+Beds+Guests.Included+Accommodates
+Features+Number.of.Reviews)
summary(model_amarquito_s)
```

```
{r}  
step(object = model_amarito_s, direction = "both")
```



Fórmula

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_k X_k$$

DATA W MEN'S COMPANY

Selección	NOMBRE / CANTIDAD DE LISTINGS		REGRESIÓN LINEAL				
	Base		Predictores aplicados	Fórmula $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$	Multiple R ^{2*} - p-value	Precio apartamento	
	Nombre	Total				Host	Superhost
Dataset	airbnb-listings.csv	14780	-	-	-	-	-
Madrid	df_madrid	6970	Bedrooms + Beds + Accommodates + Features + Number.of.Reviews	$Y = 30.59119 + 17.32385 * Br + 3.96598 * B + 6.30426 * A + 6.22138 * F - 0.15766 * NR$	0.3389	\$111,00	\$117,22
					< 2.2e-16		
Host	df_host	6068	Bedrooms + Beds + Guests.Included + Accommodates + Number.of.Reviews	$Y = 31.55054 + 17.24586 * Br + 4.24305 * B - 1.18408 * GI + 6.58241 * A - 0.17421 * NR$	0.3306	\$111,65	-
					< 2.2e-16		
Superhost	df_superhost	902	Bedrooms + Guests.Included + Accommodates + Number.of.Reviews	$Y = 30.67782 + 17.98233 * Br + 5.59964 * GI + 6.02448 * A - 0.07431 * NR$	0.4383	-	\$108,39
					< 2.2e-16		
Dº Centro (c/ de Amparo)	df_amparo	4311	Bedrooms + Beds + Guests.Included + Accommodates + Features + Number.of.Reviews	$Y = 32.3670 + 17.8795 * Br + 3.3104 * B + 1.3557 * GI + 6.1332 * A + 6.2986 * F - 0.1518 * NR$	0.4078	\$112,90	\$119,20
Bº Embajadores (c/ de Amparo)	df_amparito	1099	Bedrooms + Beds + Accommodates + Features + Number.of.Reviews	$Y = 29.1843 + 14.1807 * Br + 2.5654 * B + 6.0246 * A + 10.7265 * F - 0.1191 * NR$	0.4636	\$98,82	\$109,55
Dº Arganzuela (c/ de Méndez Álvaro)	df_alvaro	303	Bedrooms + Guests.Included + Accommodates + Number.of.Reviews	$Y = 35.2937 + 9.3598 * Br - 2.9502 * GI + 6.9069 * A - 0.1186 * NR$	0.2925	\$92,50	\$92,50
					< 2.2e-16		
Bº Palos de Moguer (c/ de Méndez Álvaro)	df_alvarito	114	Bedrooms + Guests.Included + Accommodates + Number.of.Reviews	$Y = 35.8010 + 10.6005 * Br - 8.5652 * GI + 7.8809 * A - 0.2023 * NR$	0.3888	\$95,75	\$95,75
					4.947e-11		

*En regresiones lineales múltiples, suele ser difícil conseguir un coeficiente de determinación múltiple mayor de un 30%.

DATA W MEN'S COMPANY

Selección	Nombre / CANTIDAD DE LISTINGS		REGRESIÓN LINEAL - SIN OUTLIERS				
	Limpio / Sin outliers		Predictores aplicados	Fórmula $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k + \epsilon$	Multiple R ^{2*} - p-value	Precio apartamento	
	Nombre	Total				Host	Superhost
Dataset	<i>airbnb-listings_clean_R.csv</i>	13264	-	-	-	-	
Madrid	<i>madrid_s</i>	5229	Bedrooms + Beds + Guests.Included + Accommodates + Features + Number.of.Reviews	$Y = 49.34626 + 7.18268 * Br + 1.77273 * B - 1.73606 * Gl + 4.91038 * A + 8.9324 * F - 0.23564 * NR$	0.2012 ≤ 2.2e-16	\$94,98	\$103,92
Host	<i>host_s</i>	4497	Bedrooms + Beds + Guests.Included + Accommodates + Number.of.Reviews	$Y = 48.66326 + 7.25955 * Br + 2.06097 * B - 1.11485 * Gl + 4.85718 * A - 0.29265 * NR$	0.2058 ≤ 2.2e-16	\$95,33	
Superhost	<i>super_s</i>	644	Bedrooms + Guests.Included + Accommodates	$Y = 57.976 + 8.573 * Br + 4.891 * Gl + 4.984 * A$	0.1562 ≤ 2.2e-16		\$100,14
Dº Centro (c/ de Amparo)	<i>amparo_s</i>	3495	Bedrooms + Beds + Guests.Included + Accommodates + Features + Number.of.Reviews	$Y = 45.2834 + 9.4199 * Br + 1.4329 * B + 2.2162 * Gl + 4.1568 * A + 7.0679 * F - 0.1714 * NR$	0.305 ≤ 2.2e-16	\$96,15	\$101,21
Bº Embajadores (c/ de Amparo)	<i>amparito_s</i>	922	Bedrooms + Beds + Guests.Included + Accommodates + Features + Number.of.Reviews	$Y = 38.4211 + 9.0340 * Br + 2.5862 * B + 1.2239 * Gl + 3.6121 * A + 9.7052 * F - 0.1477 * NR$	0.3544 ≤ 2.2e-16	\$84,56	\$94,26
Dº Arganzuela (c/ de Méndez Álvaro)	<i>alvaro_s</i>	219	Bedrooms + Beds + Features + Number.of.Reviews	$Y = 46.1461 + 4.1607 * Br + 4.6348 * B + 6.3421 * F - 0.1417 * NR$	0.1387 1.666e-05	\$63,74	\$70,08
Bº Palos de Moguer (c/ de Méndez Álvaro)	<i>alvarito_s</i>	86	Bedrooms + Beds+ Guests.Included + Accommodates	$Y = 37.900 + 7.713 * Br + 5.142 * B - 7.518 * Gl + 4.947 * A$	0.3454 5.275e-07	\$85,77	\$85,77

*En regresiones lineales múltiples, suele ser difícil conseguir un coeficiente de determinación múltiple mayor de un 30%.

7. RESULTADOS

Selección	REGRESIÓN LINEAL		REGRESIÓN LINEAL - SIN OUTLIERS		TABLEAU	
	Precio apartamento		Precio apartamento		Precio apartamento	
	Host	Superhost	Host	Superhost	Host	Superhost
Bº Embajadores (c/ de Amparo)	\$98,82	\$109,55	\$84,56	\$94,26	\$80	\$95,00
Bº Palos de Moguer (c/ de Méndez Álvaro)	\$95,75	\$95,75	\$85,77	\$85,77	\$74,50	\$95,00

REVIEWS:

- Barrio Embajadores:

$$Y=29.1843+14.1807*BR+2.5654*B+6.0246*A+10.7265*F-0.1191*NR$$

- Barrio Palos de Moguer:

$$Y=35.8010+10.6005*BR-8.5652*GI+7.8809*A-0.2023*NR$$

8. CONCLUSIONES

Precios base

- Barrio Embajadores:
 - Host: **\$80,00 - 84,16**
 - Superhost: **\$94,26 - 95,00**
- Barrio Palos de Moguer:
 - Host: **\$74,90 - 85,77**
 - Superhost: **\$95,00 - \$95,75**
(sin regresión lineal)
 - El precio es similar en ambos barrios.
 - La regresión lineal de ambos y Tableau indica:
(-) precio~nº de reviews
 - Ser superhost implica
(+) ~ \$15,00

¿QUÉ TE HA APORTADO EL DESARROLLAR ESTE PROYECTO?

- Satisfacción de superar retos y ser resolutivas
- Consolidación de conocimientos mediante práctica
- Trabajo en equipo implementando SCRUM
- Hacernos amigas de la frustración,



¿QUÉ HAS APRENDIDO?

- Saber cómo y en qué momento aplicar las herramientas
- A trabajar en un proyecto pseudorreal
- A conocer los límites personales
- A implementar nuevas funciones de herramientas

¿QUÉ ES LO QUE NO VOLVERÍAS A HACER DE LA MISMA MANERA?

- Big Data es aprendizaje diario, tenemos mucho que descubrir y aprender
HARÍAMOS IGUAL PERO CON MÁS TIEMPO

¿QUÉ COSAS SEGUIRÍAS HACIENDO EN EL FUTURO PARA MEJORAR EL PROYECTO?

- Continuar con el modelado realizando regresiones logísticas.
- Conseguir más métricas continuas o métricas temporales para mejorar los resultados

¡SEGUIR TRABAJANDO Y MEJORANDO NUESTRAS SKILLS!

DATA WOMEN'S COMPANY

Enlaces de interés:



Fórmula
 $Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_k \lambda$



DATA WOMEN'S COMPANY

WHERE DATA MAKES SENSE

Nos gustaría dar las gracias a todo el equipo de KeepCoding, Glovo y a todas las compañeras que nos han acompañado en este increíble viaje.

CON CARIÑO

DATA WOMEN'S COMPANY
WHERE DATA MAKES SENSE