





Interpretación Técnica de la Matriz de Correlación (Sevilla):

La matriz de correlación revela las relaciones lineales entre variables clave en los alojamientos de Airbnb en Sevilla. A continuación, se desglosan los hallazgos principales:

1. Relación Precio-Valoración

- **Correlación: -0.01**
  - **Interpretación:** Al igual que en Madrid, no existe una relación lineal significativa entre precio y valoración.
  - **Implicación:** Los precios altos no garantizan mejores valoraciones, y las valoraciones altas no están vinculadas a precios elevados.
  - **Posible causa:** Factores no analizados (ej: calidad del servicio, ubicación exacta, limpieza) podrían influir más en las valoraciones.

2. Precio vs Tipo de Habitación

- **Entire home/apt (Correlación: 0.06)**
  - Las propiedades completas tienen una correlación positiva muy débil con el precio, similar a Madrid pero menos pronunciada.
- **Private room (Correlación: -0.06)**
  - Ligera tendencia a precios más bajos en habitaciones privadas.
- **Relación inversa (Entire home/apt vs Private room: -0.98)**
  - Interpretación: Correlación negativa casi perfecta, confirmando la exclusividad mutua de las categorías (un alojamiento no puede ser ambas cosas).

3. Dormitorios vs Precio

- **Correlación: 0.10**
  - **Hallazgo:** Igual que en Madrid, el número de dormitorios tiene un impacto mínimo en el precio.
  - **Insight:** En Sevilla, el tamaño de la propiedad no es un driver principal de precio. La ubicación o amenities podrían ser más relevantes.

4. Valoración vs Otras Variables

- **Todas las correlaciones cercanas a 0** (entre -0.04 y 0.07).
  - **Conclusión:** La valoración no está ligada linealmente a:
    - Precio.
    - Dormitorios.
    - Tipo de habitación.
  - **Recomendación:** Analizar variables cualitativas (ej: comentarios de reseñas) para identificar drivers ocultos.

5. Relaciones entre Tipos de Habitación

- **Entire home/apt vs Private room: -0.98**
  - **Validación técnica:** Coherencia con la codificación dummy, donde un alojamiento solo puede pertenecer a una categoría.

Diferencias Clave vs Madrid

- Dormitorios y Entire home/apt (Correlación: 0.20 en Sevilla vs 0.18 en Madrid):**
  - En Sevilla, las propiedades completas tienden ligeramente a tener más dormitorios, pero la relación sigue siendo débil.
- Valoración y Dormitorios (Correlación: 0.07 en Sevilla vs -0.01 en Madrid):**
  - Aunque insignificante, sugiere que en Sevilla podría haber una mínima tendencia a valorar mejor propiedades con más dormitorios.

Conclusiones para el Informe:

- Estrategia de Precios:** En Sevilla, el tipo de habitación influye menos en el precio que en Madrid. Las propiedades completas no justifican precios significativamente más altos.
  - Valoraciones:** Al igual que en Madrid, se requieren análisis no lineales o cualitativos para entender los motivos detrás de las puntuaciones.
  - Segmentación de mercado:** La elección entre "propiedad completa" y "habitación privada" es igualmente excluyente en ambas ciudades, pero en Sevilla la correlación precio-tipo es aún más tenue.
- Recomendación final:** Realizar un estudio comparativo entre ciudades para identificar si las diferencias en correlaciones reflejan variaciones culturales o de mercado.

Este análisis proporciona una base cuantitativa para optimizar estrategias comerciales en Sevilla, destacando la necesidad de enfoques complementarios para entender las valoraciones.

Visualizaciones

Gráfico Comparativo

Distribución del precio por tipo de alojamiento y ciudad

```
In [ ]: # Filtra los precios que son menores a 1000 euros
airbnb_variables_objetivos_limpio_filtrado = airbnb_variables_objetivos_limpio[airbnb_variables_obj

# Gráfico comparativo de distribución del precio por alojamiento y ciudad
sns.boxplot(data=airbnb_variables_objetivos_limpio_filtrado, x='room_type', y='price', hue='city')
plt.title = "Distribución del precio por tipo de alojamiento y ciudad"
plt.xticks(rotations=45)
plt.savefig('boxplot_filtrado.png')
plt.show()
```

Estadística Robusta (Pope)

Aunque ya se aplicó una limpieza de precios mediante el rango intercuartílico (IQR), se ha decidido conservar esta gráfica como herramienta **visual de validación**. Su objetivo es comparar medidas robustas (mediana e IQR) frente a medidas sensibles a valores extremos (media y desviación estándar) para verificar si los precios altos que persisten son esperables, especialmente en ciudades como Madrid o Sevilla.

No se utilizará esta gráfica como criterio adicional de filtrado, sino como **respaldo visual y exploratorio** del comportamiento de los precios por ciudad.

```
In [ ]: # 1. Filtrar los precios menores a 1000
airbnb_variables_objetivos_limpio_filtrado = airbnb_variables_objetivos_limpio[airbnb_variables_obj

# 2. Agrupa por ciudad y calcula estadísticas básicas
stats = airbnb_variables_objetivos_limpio_filtrado.groupby('city')['price'].agg(['mean', 'std', 'me

# 3. Calcula los cuartiles y el rango intercuartílico (IQR)
q1 = airbnb_variables_objetivos_limpio_filtrado.groupby('city')['price'].quantile(0.25)
q3 = airbnb_variables_objetivos_limpio_filtrado.groupby('city')['price'].quantile(0.75)
stats['iqr'] = q3 - q1

# 4. Prepara las propiedades para las barras
cities = stats.index.tolist()
x = range(len(cities))
width = 0.35

# 5. Dibuja barras con barras de error
fig, ax = plt.subplots(figsize=(8, 5))
ax.bar([1 + width/2 for i in x],
       stats['mean'],
       width,
       yerr=stats['std'],
       capsize=5,
       label='Media ± Desviación estándar',
       color='skyblue')
ax.bar([1 + width/2 for i in x],
       stats['median'],
       width,
       yerr=stats['iqr'],
       capsize=5,
       label='Mediana ± IQR',
       color='lightgreen')

# 6. Ajustes Finales
ax.set_xticks(x)
ax.set_xticklabels(cities, rotation=45)
ax.set_title('Estadísticas Robusta vs No Robusta (Precio por ciudad)')
ax.set_ylabel('Precio (€)')
ax.legend()
plt.tight_layout()
plt.savefig('EstadisticaRVSnm.png')
plt.show()
```

Análisis de Correlación entre Variables de Alojamiento en Sevilla y Madrid

```
In [ ]: # Filtrar los datos para Sevilla y Madrid
sevilla_df = airbnb_variables_objetivos_limpio[airbnb_variables_objetivos_limpio['city'] == 'Sevilla']
madrid_df = airbnb_variables_objetivos_limpio[airbnb_variables_objetivos_limpio['city'] == 'Madrid']

# Crear matriz de correlación para Sevilla
corr_sevilla = sevilla_df[['price', 'bedrooms', 'valoracion']].corr()

# Crear matriz de correlación para Madrid
corr_madrid = madrid_df[['price', 'bedrooms', 'valoracion']].corr()

# Crear una figura con dos subgráficas (2 filas, 1 columna)
plt.figure(figsize=(16, 6))

# Primer gráfico: Correlación Sevilla
plt.subplot(1, 2, 1) # 1 fila, 2 columnas, 1er gráfico
sns.heatmap(corr_sevilla,
            annot=True,
            cmap='coolwarm',
            vmin=-1,
            vmax=1)

plt.gca().set_title("Correlación entre variables - Sevilla")

# Segundo gráfico: Correlación Madrid
plt.subplot(1, 2, 2) # 1 fila, 2 columnas, 2do gráfico
sns.heatmap(corr_madrid,
            annot=True,
            cmap='coolwarm',
            vmin=-1,
            vmax=1)

plt.gca().set_title("Correlación entre variables - Madrid")

# Ajustar el espacio entre los gráficos
plt.tight_layout()

# Mostrar las gráficas
plt.show()

#### Identificar pares de variables con alta correlación (umbral > 0.7 o < -0.7)

# Filtrar pares con alta correlación (umbral > 0.7 o < -0.7)
# Madrid
correlated_pairs_madrid = corr_madrid.abs().unstack().sort_values(ascending=False)
high_correlation_pairs_madrid = correlated_pairs_madrid[(correlated_pairs_madrid > 0.7) & (correlat
print(high_correlation_pairs_madrid)

#Sevilla
correlated_pairs_sevilla = corr_sevilla.abs().unstack().sort_values(ascending=False)
high_correlation_pairs_sevilla = correlated_pairs_sevilla[(correlated_pairs_sevilla > 0.7) & (corre
print(high_correlation_pairs_sevilla)
```

Visualizar la correlación de variables clave en el dataset (airbnb\_variables\_objetivos\_limpio)

```
In [ ]: # Eliminar variable title si existe
if 'title' in globals():
    del title

# =====
# ANALISIS PARA MADRID
# =====
data_madrid = airbnb_variables_objetivos_limpio[airbnb_variables_objetivos_limpio['city'] == 'Madr

# Configuración de estilo
sns.set_theme(style="whitegrid", palette="pastel")
fig, axes = plt.subplots(2, 2, figsize=(16, 12))
fig.suptitle('Análisis de Propiedades en Madrid', fontsize=18, fontweight='bold', color='#2C3E50')

# 1. Distribución de Precios por Tipo de Habitación
sns.histplot(data=data_madrid,
             x='price',
             hue='room_type',
             bins=30,
             ax=axes[0, 0])
axes[0, 0].set_title('¿Cómo se distribuyen los precios entre tipos de habitación?\n¿Existen diferenc
             fontsize=12, pad=10)
axes[0, 0].set_xlabel('Precio (€)', fontsize=10)
axes[0, 0].set_ylabel('Frecuencia', fontsize=10)
axes[0, 0].set_xlim(0, 300)

# 2. Relación Precio-Valoración
sns.scatterplot(data=data_madrid,
               x='valoracion',
               y='price',
               hue='room_type',
               alpha=0.7,
               ax=axes[0, 1])
axes[0, 1].set_title('¿Las propiedades más caras tienen mejor valoración?\n¿Varia esta relación por
             fontsize=12, pad=10)
axes[0, 1].set_xlabel('Valoración (1-5)', fontsize=10)
axes[0, 1].set_ylabel('Precio (€)', fontsize=10)

# 3. Distribución de Dormitorios
sns.histplot(data=data_madrid,
             x='bedrooms',
             bins=15,
             kde=True,
             ax=axes[1, 0])
axes[1, 0].set_title('¿El tamaño de las propiedades es el mismo?\n¿Cómo se distribuyen los dor
             fontsize=12, pad=10)
axes[1, 0].set_xlabel('Dormitorios', fontsize=10)
axes[1, 0].set_ylabel('Frecuencia', fontsize=10)
axes[1, 0].set_xlim(range(0, 11))

# 4. Precios en Barrios Top
top_barrios_madrid = data_madrid['neighbourhood'].value_counts().index[:5]
sns.boxplot(data=data_madrid[data_madrid['neighbourhood'].isin(top_barrios_madrid)],
            x='neighbourhood',
            y='price',
            ax=axes[1, 1])
axes[1, 1].set_title('¿Qué barrios tienen precios más estables?\n¿Dónde encontramos outliers signifi
             fontsize=12, pad=10)
axes[1, 1].set_xlabel('Precio (€)', fontsize=10)
axes[1, 1].tick_params(axis='x', rotation=45)

plt.tight_layout(rect=[0, 0.03, 1, 0.95])
plt.show()

# =====
# ANALISIS PARA SEVILLA
# =====
data_sevilla = airbnb_variables_objetivos_limpio[airbnb_variables_objetivos_limpio['city'] == 'Sevi

fig, axes = plt.subplots(2, 2, figsize=(16, 12))
fig.suptitle('Análisis de Propiedades en Sevilla', fontsize=18, fontweight='bold', color='#2C3E50')

# 1. Distribución de Precios por Tipo de Habitación (Sevilla)
sns.histplot(data=data_sevilla,
             x='price',
             hue='room_type',
             bins=30,
             ax=axes[0, 0])
axes[0, 0].set_title('¿Cómo comparan los precios con Madrid?\n¿Qué tipo de habitación domina el merc
             fontsize=12, pad=10)
axes[0, 0].set_xlabel('Precio (€)', fontsize=10)
axes[0, 0].set_ylabel('Frecuencia', fontsize=10)
axes[0, 0].set_xlim(0, 300)

# 2. Relación Precio-Valoración (Sevilla)
sns.scatterplot(data=data_sevilla,
               x='valoracion',
               y='price',
               hue='room_type',
               alpha=0.7,
               ax=axes[0, 1])
axes[0, 1].set_title('¿Cómo mantiene la relación precio-valoración?\n¿Hay diferencias con Madrid?',
             fontsize=12, pad=10)
axes[0, 1].set_xlabel('Valoración (1-5)', fontsize=10)
axes[0, 1].set_ylabel('Precio (€)', fontsize=10)

# 3. Distribución de Dormitorios (Sevilla)
sns.histplot(data=data_sevilla,
             x='bedrooms',
             bins=15,
             kde=True,
             ax=axes[1, 0])
axes[1, 0].set_title('¿El tamaño de las propiedades es diferente?\n¿Qué distribución de dormitorios
             fontsize=12, pad=10)
axes[1, 0].set_xlabel('Dormitorios', fontsize=10)
axes[1, 0].set_ylabel('Frecuencia', fontsize=10)
axes[1, 0].set_xlim(range(0, 11))

# 4. Precios en Barrios Top (Sevilla)
top_barrios_sevilla = data_sevilla['neighbourhood'].value_counts().index[:5]
sns.boxplot(data=data_sevilla[data_sevilla['neighbourhood'].isin(top_barrios_sevilla)],
            x='neighbourhood',
            y='price',
            ax=axes[1, 1])
axes[1, 1].set_title('¿Cómo se comparan los barrios populares con Madrid?\n¿Patrones de precios dife
             fontsize=12, pad=10)
axes[1, 1].set_xlabel('Precio (€)', fontsize=10)
axes[1, 1].tick_params(axis='x', rotation=45)

plt.tight_layout(rect=[0, 0.03, 1, 0.95])
plt.show()
```

precios medios por tipo de habitación y ciudad

Interpretación potencial:

- ¿Los apartamentos completos son significativamente más caros en Madrid?
- ¿Qué ciudad ofrece mejores precios en habitaciones compartidas?
- ¿Existe mayor variación de precios en algún tipo específico de alojamiento?

```
In [ ]: # PASO 1: LIMPIEZA DE VARIABLES (EJECUTAR ESTO PRIMERO)
# =====
import matplotlib.pyplot as plt

# Verificar y eliminar variables conflictivas
conflict_vars = ['title', 'Title', 'TITLE']
for var in conflict_vars:
    if var in globals():
        del globals()[var]
        print(f"Variable conflictiva '{var}' eliminada")

# Reiniciar configuración de matplotlib
plt.close('all')
plt.rcParams.update(plt.rcParamsDefault)

# PASO 2: CÓDIGO GRÁFICO CORREGIDO (EJECUTAR DESPUÉS DEL PASO 1)
# =====
import seaborn as sns

# Calcular precios medios
grouped = airbnb_variables_objetivos_limpio.groupby(['room_type', 'city'])['price'].mean().reset_in

# Configurar figura
fig, ax = plt.subplots(figsize=(12, 6))
sns.set_theme(style="whitegrid", palette="pastel")

# Crear gráfico
ax = sns.barplot(
    x='room_type',
    y='price',
    data=grouped,
    errorbar=None,
    palette=("Madrid": "#e74c3c", "Sevilla": "#2980b9"),
    saturation=0.85 # Intensidad de color aumentada
)

# Personalización avanzada
ax.set_title('Comparativa de precios medios\nMadrid vs Sevilla',
            fontsize=16,
            pad=20,
            fontweight='bold',
            color='#2c3e50',
            loc='left' # Alineación izquierda
)

ax.set_xlabel('Tipo de alojamiento',
            fontsize=12,
            labelpad=15,
            fontStyle='italic'
)

ax.set_ylabel('Precio medio ($)',
            fontsize=12,
            labelpad=15,
            fontStyle='italic'
)

# Rotación y formato de ejes
plt.xticks(
    rotation=20,
    fontsize=11,
    fontfamily='DejaVu Sans' # Fuente más legible
)

# Etiquetas de valor mejoradas
for p in ax.patches:
    ax.text(
        p.get_x() + p.get_width() / 2.,
        p.get_height() + 5, # Posición vertical ajustada
        f"${p.get_height():.0f}",
        ha='center',
        va='center',
        fontsize=10,
        color='#34495e'
    )

# Leyenda personalizada
leg = ax.legend(
    title='Ciudad',
    title_fontsize='13',
    fontsize=11,
    frameon=True,
    shadow=True,
    facecolor='#f9f9fa' # Fondo claro
)

# Ajustes finales
plt.tight_layout()
plt.show()
```

Conversión del Dataset a Formato JSON

¿Cómo validar que los archivos .json tienen estructura válida? Para asegurarte de que los archivos .json estén bien estructurados y sean compatibles con MongoDB, debes cumplir con lo siguiente:

Estructura válida JSON Puedes validarlos en línea en herramientas como:

<https://jsonlint.com/>

<https://jsonmatter.org/>

Compatible con MongoDB MongoDB puede importar:

- Un JSON por línea (ndjson o JSONL)
- O un archivo con un array de objetos JSON, así:

```
{
  { "nombre": "Juan", "edad": 35 },
  { "nombre": "Ana", "edad": 29 }
}
```

Recomendación: Usa el formato de array de objetos si vas a cargar con mongolimport.

```
In [ ]: # Información del dataset limpio
airbnb_variables_objetivos_limpio.info()

In [ ]: # Limitar la cantidad de filas a 10 filas para evitar que los archivos .json sean demasiado grandes
airbnb_madrid_y_sevilla_limited = airbnb_variables_objetivos_limpio.head(10)

# Dataframe airbnb_variables_objetivos_limpio completo
airbnb_madrid_y_sevilla_completo = airbnb_variables_objetivos_limpio.copy()

In [ ]: # Enviar todo el contenido del dataframe a un archivo .json (No se limita la cantidad de filas)
#alojamientos,turisticos.to_json('airbnb_variables_objetivos_limpio.json', orient='records', force_

airbnb_madrid_y_sevilla_limited.to_json('airbnb_madrid_y_sevilla_limited.json', orient='records',
airbnb_madrid_y_sevilla_completo.to_json('airbnb_madrid_y_sevilla_completo.json', orient='records'
```

Exportar cuaderno a formato HTML y a formato PDF

Descarga el archivo generado directamente desde Colab

```
In [ ]: import shutil
from google.colab import files

# Función para limpiar y validar el notebook
def fix_notebook_metadata(file_path):
    with open(file_path, 'r', encoding='utf-8') as file:
        notebook = json.load(file)

        for cell in notebook.get('cells', []):
            if 'outputs' in cell:
                for output in cell['outputs']:
                    # Agregar 'metadata' si no existe
                    if 'metadata' not in output:
                        output['metadata'] = {}

        # Guardar el notebook limpio
        with open(file_path, 'w', encoding='utf-8') as file:
            json.dump(notebook, file, indent=2, ensure_ascii=False)

# Rutas de archivos
notebook_path = "%content/drive/MyDrive/Cursos/UNIR - BigDataAI/proyectos/bases-de-datos-para-e-Big
html_output_path = "%content/drive/MyDrive/Cursos/UNIR - BigDataAI/proyectos/bases-de-datos-para-e-
pdf_output_path_viahtml = "%content/drive/MyDrive/Cursos/UNIR - BigDataAI/proyectos/bases-de-datos-

# Corregir el notebook
print("\nCorrigiendo y validando el notebook...")
fix_notebook_metadata(notebook_path)

# Exportar a HTML
print("\nExportando a formato HTML...")
!jupyter nconvert --to html --output "html_output_path" "Notebook.path"

# Exportar a PDF usando PDFVIAHTML
print("\nExportando a formato PDF usando nconvert con PDFVIAHTML...")
!jupyter nconvert --to pdfviahtml --output "pdf_output_path_viahtml" "Notebook.path"

# Descargar archivos generados
print("\nDescargando los archivos generados...")
try:
    files.download(html_output_path) # Descargar HTML
    files.download(pdf_output_path_viahtml) # Descargar PDF generado con PDFVIAHTML
except Exception as e:
    print(f"Error descargando los archivos: {e}")

print("\nExportación y descarga completadas.")
```